



Swin Transformer 从零解读

主讲人：DASOU

回复【Swin】获取对应PPT



扫码关注微信公众号

NLP从入门到放弃

文章周更

知识分享

一起进步

求关注，求点赞，求一切！！

1. 回顾TRM和VIT模型，理清SwinTRM的创新点

2. SwinTRM使用到相对位置编码介绍

3. 窗口移动注意力机制SW-MSA介绍

4. Patch Merging介绍

回顾TRM模型

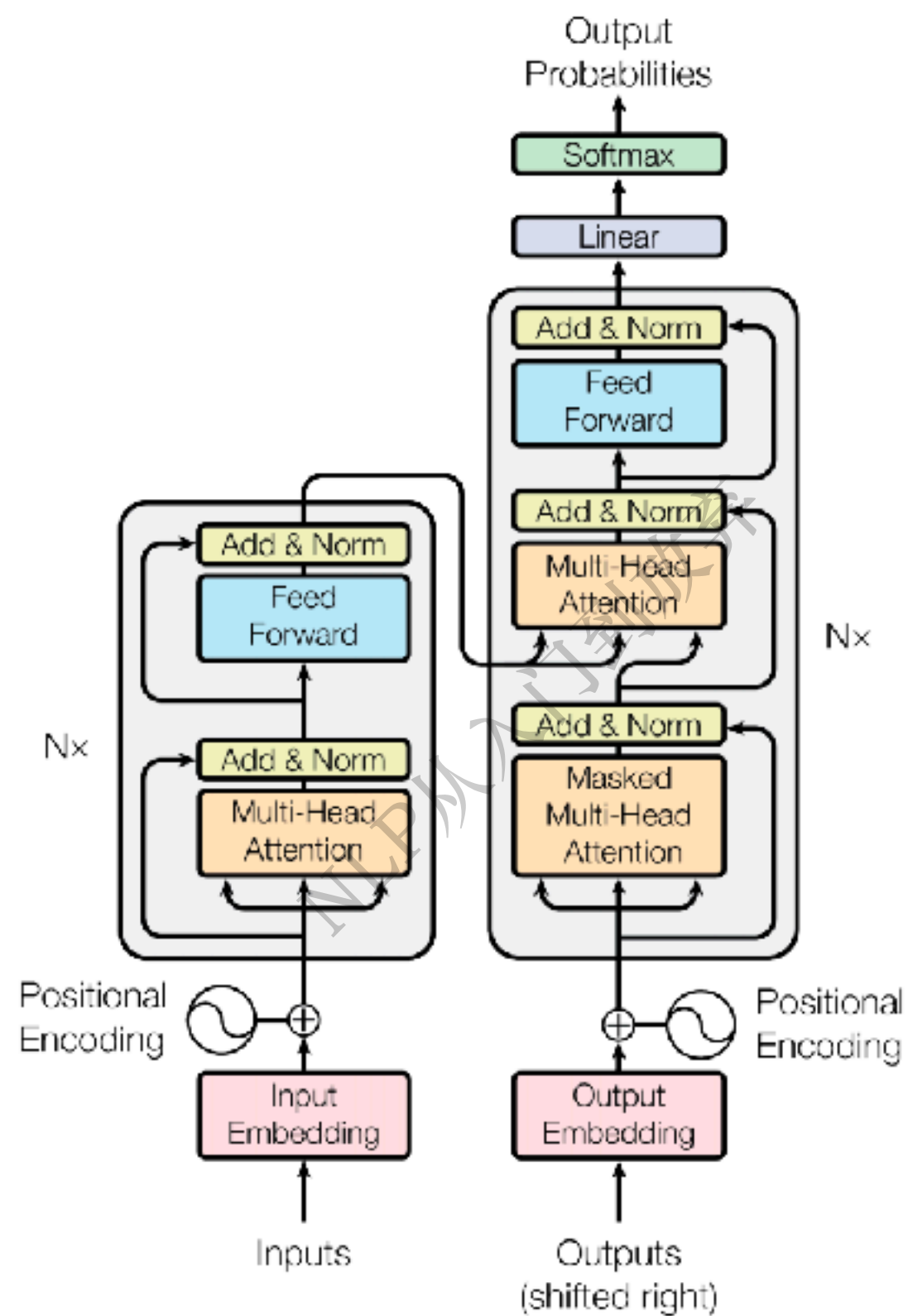
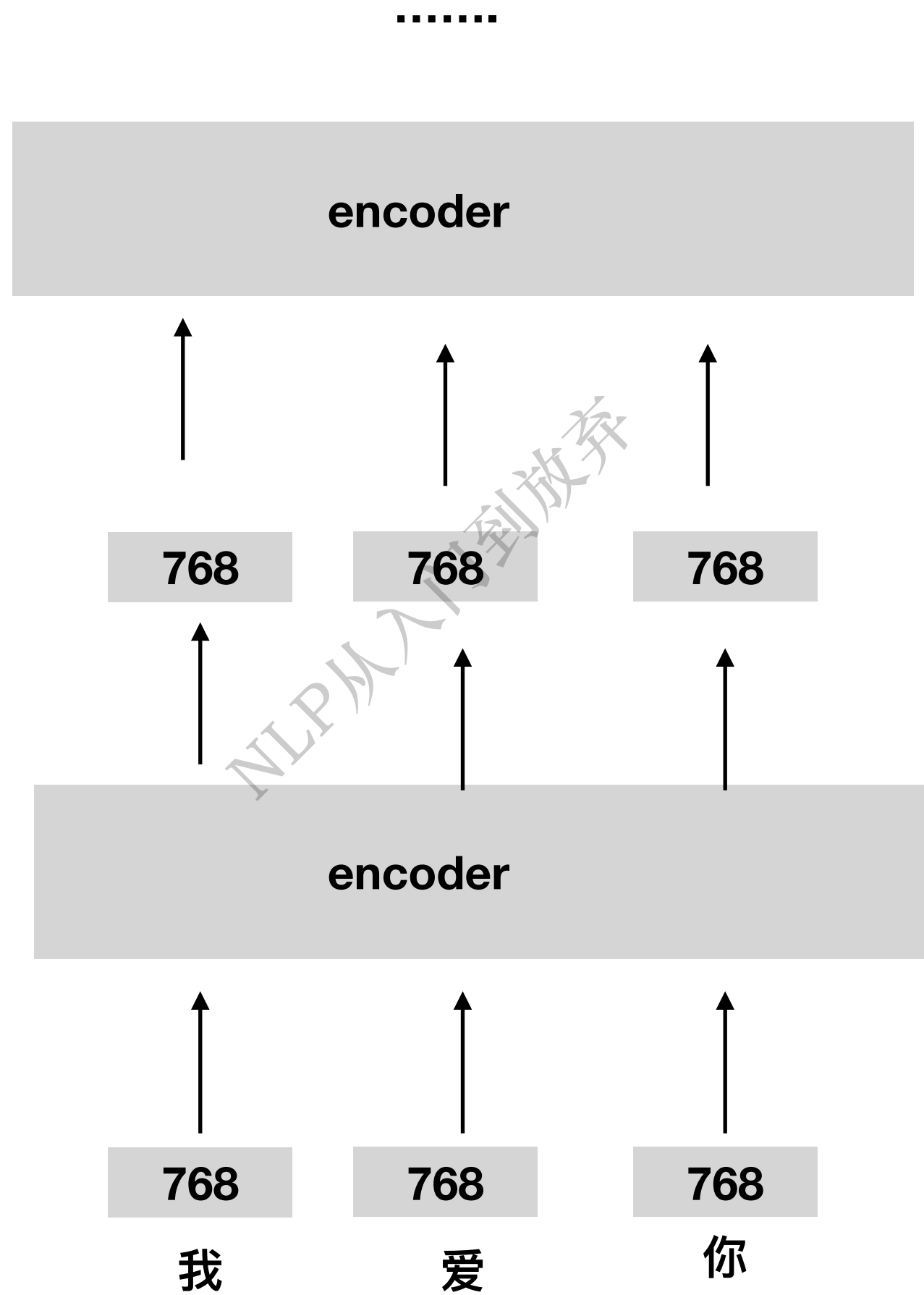
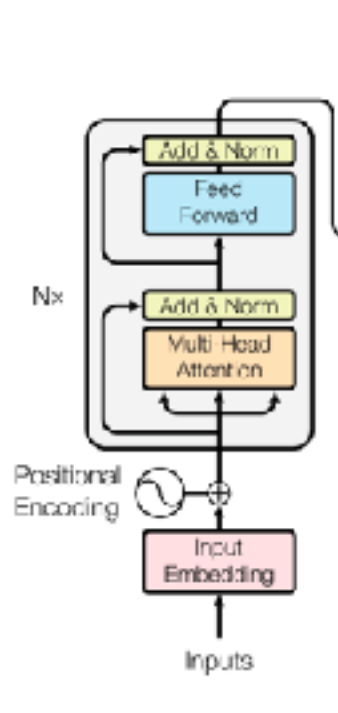


Figure 1: The Transformer - model architecture.



回顾ViT模型

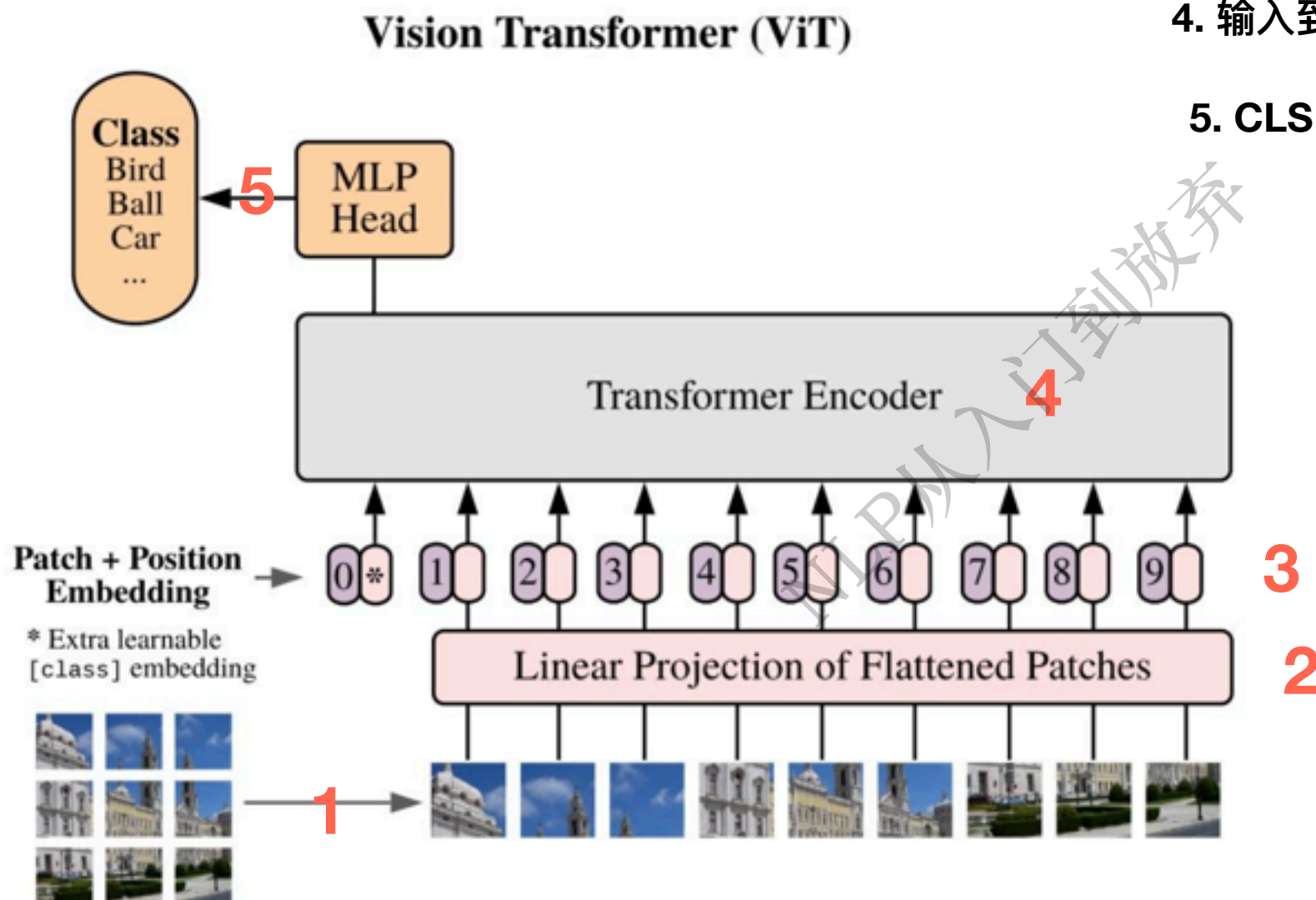
1. 图片切分为patch

2. patch转化为embedding

3.位置embedding和tokenembedding相加

4. 输入到TRM模型

5. CLS 输出做多分类任务



3.1 生成CLS符号的TokenEMB

3.2 生成所有序列的位置编码

3.3 token+位置编码

SwinTRM做到了两点

- 1. 金字塔形状：感受野是在不停的变大的**
- 2. 注意力机制放在一个窗口内部**

NLP从入门到放弃

当时我们VIT模型中怎么把图片变成token的时候聊过

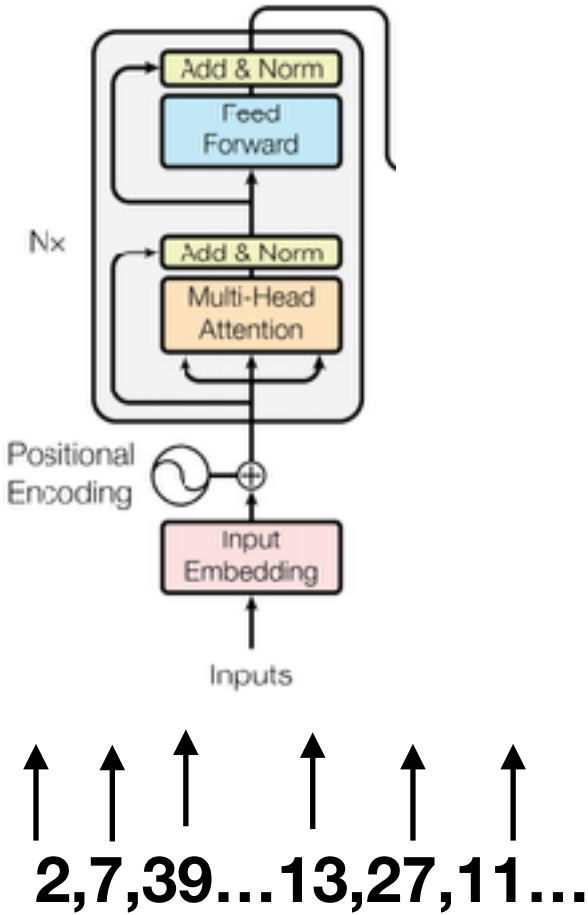
大部分人的思路



→ 224

2	7	39	
13	27	11	
39	7	2	
...			...

224



复杂度的问题

$224 \times 224 \times 1$

224

224

2	7	39	
13	27	11	
39	7	2	
...			...

→ 序列长度 = $224 \times 224 = 50176$

BERT的最大长度是512，相当于100倍

如何处理复杂度的问题？：本质上是去解决随着像素增加，复杂度平方级增长的问题；

1.局部注意力机制

有很多中方法：

2. 改进attention公式

3.....

VIT是图片化整为零：由一个像素点转为为一块像素点

一个简单的改进方式：图像化整为零，切分patch

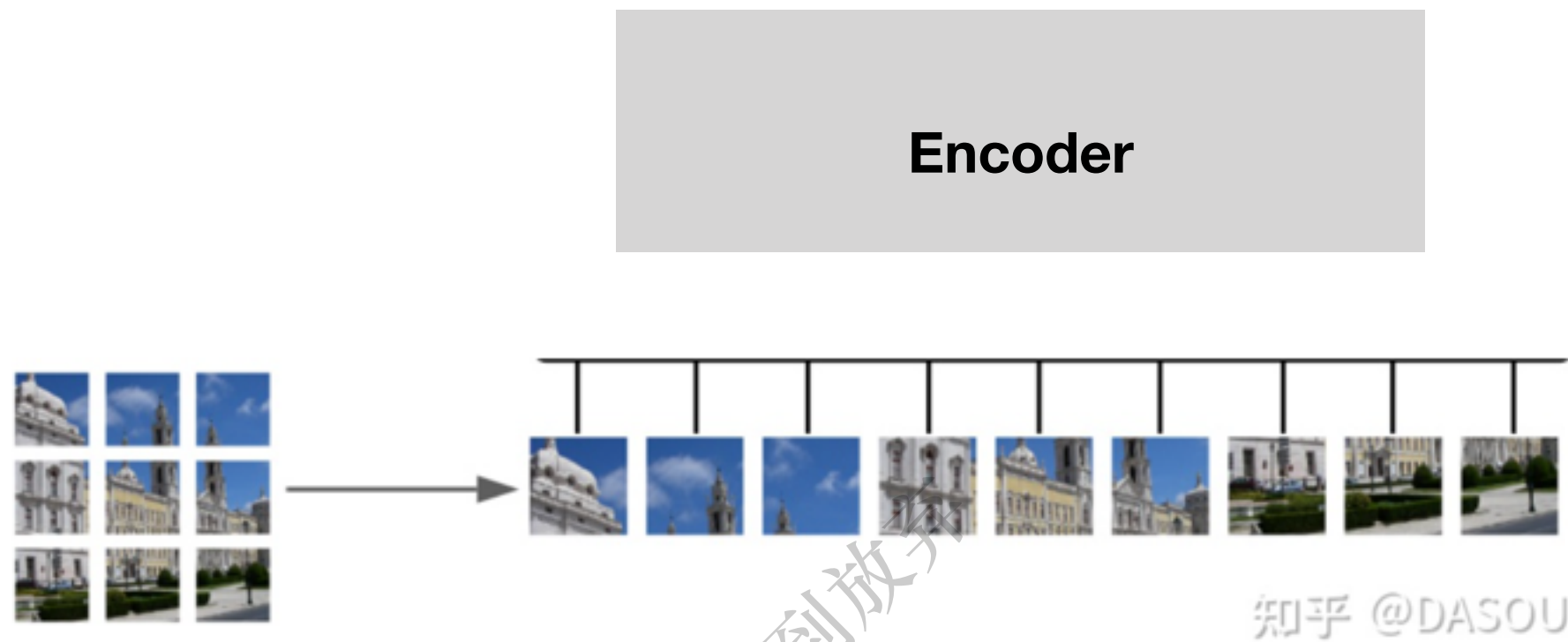
**也就说原来是一个像素点代表一个token，
现在是一大块的token一个patch作为一个token**



知乎 @DASOU

VIT的问题就是如果图像过大，patch就变大了，复杂度还是不太好

VIT



Swintrm



1: $224 \times 224 \times 3$ 到 $(224/4) \times (224/4) \times (4 \times 4 \times 3)$

来看SwinTRM的架构图

2. $(224/4) \times (224/4) \times (4 \times 4 \times 3)$ 到 $(224/4) \times (224/4) \times (96)$

3. $(224/4) \times (224/4) \times (96)$ 到 $(224/4) \times (224/4) \times (96)$

4. $(224/4) \times (224/4) \times (96)$ 到 $(224/8) \times (224/8) \times (96 \times 2 = 192)$

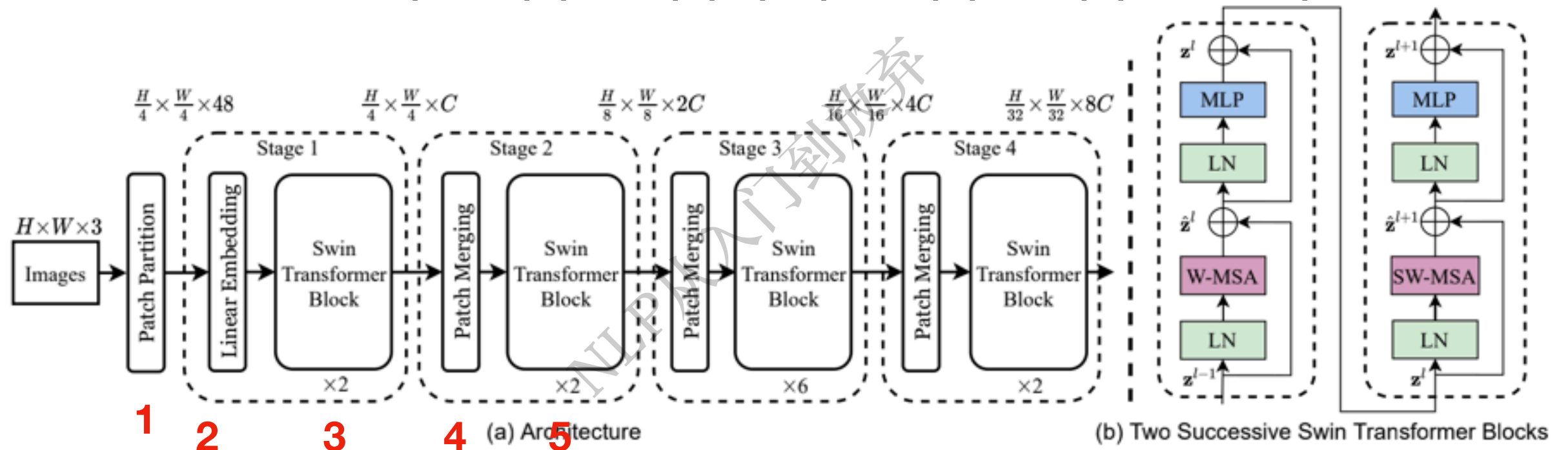
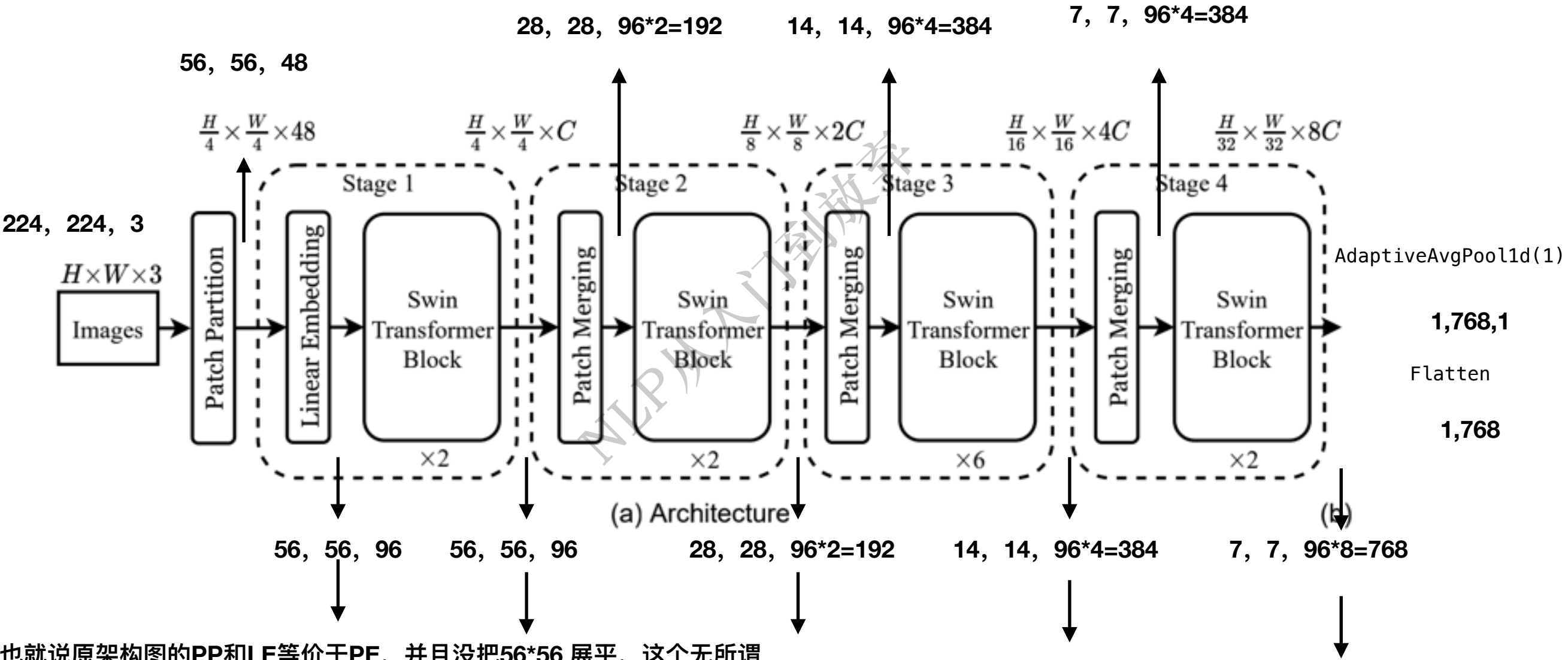


Figure 3. (a) The architecture of a Swin Transformer (Swin-T); (b) two successive Swin Transformer Blocks (notation presented with Eq. (3)). W-MSA and SW-MSA are multi-head self attention modules with regular and shifted windowing configurations, respectively.

5. $(224/8) \times (224/8) \times (96 \times 2 = 192)$ 到 $(224/8) \times (224/8) \times (96 \times 2 = 192)$

H是224, W是224



来看SwinTRM的架构图

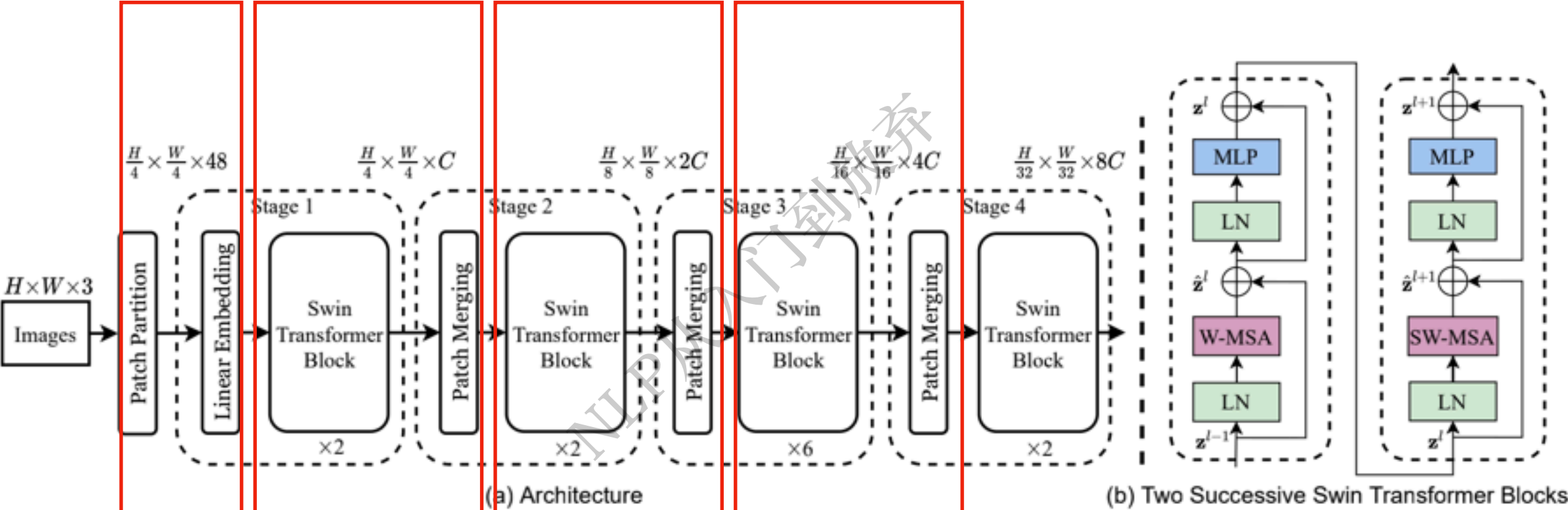
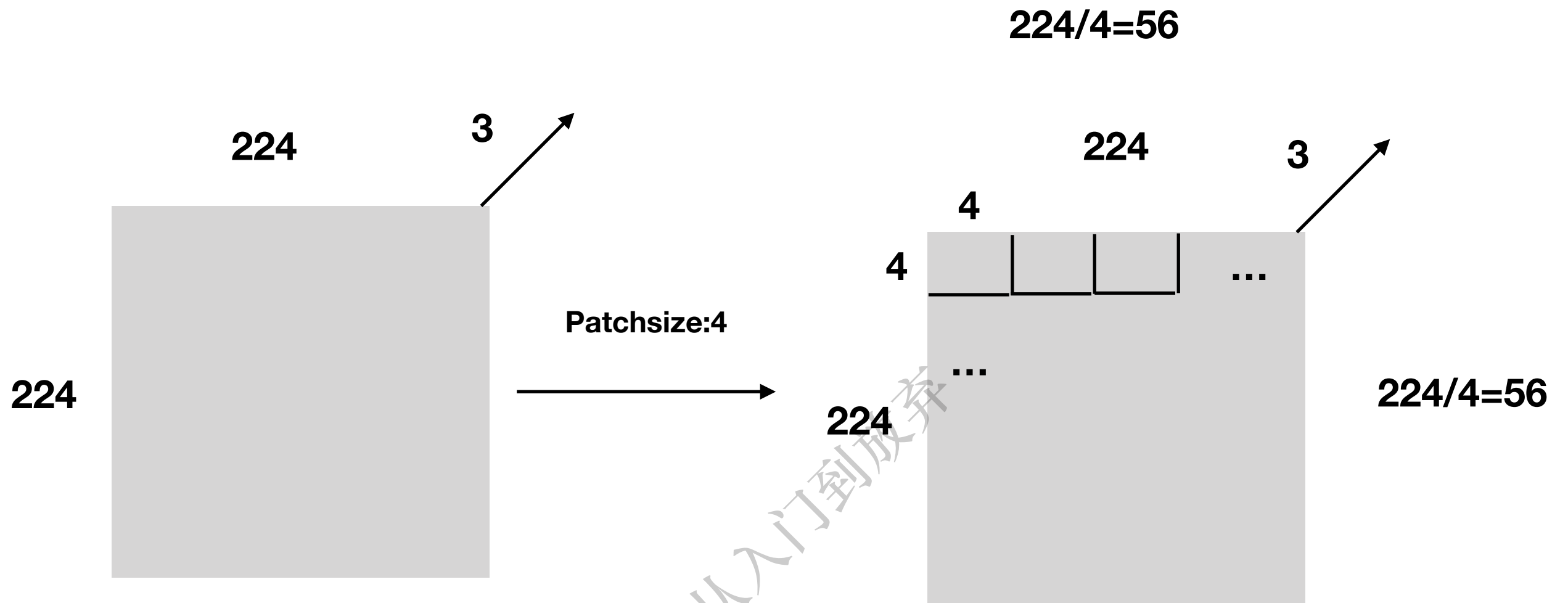


Figure 3. (a) The architecture of a Swin Transformer (Swin-T); (b) two successive Swin Transformer Blocks (notation presented with Eq. (3)). W-MSA and SW-MSA are multi-head self attention modules with regular and shifted windowing configurations, respectively.

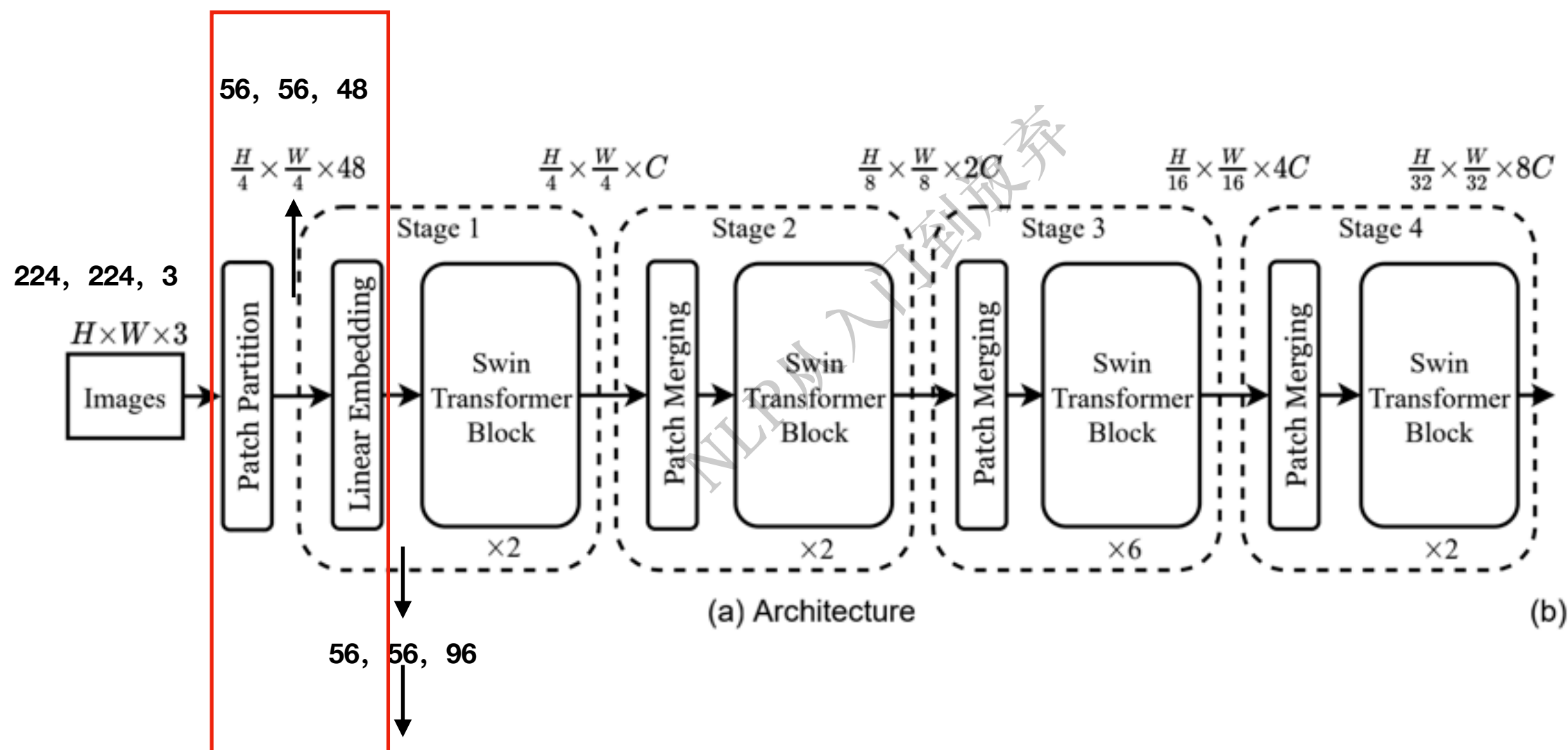


Patch数目就是： $56*56=3136$

1: $224*224*3$ 到 $(224/4)*(224/4)*(4*4*3)$

2. $(224/4)*(224/4)*(4*4*3)$ 到 $(224/4)*(224/4)*(96)$

H是224, W是224



也就是说原架构图的PP和LE等价于PE，并且没把56*56展平，这个无所谓

第二个红色框框的操作

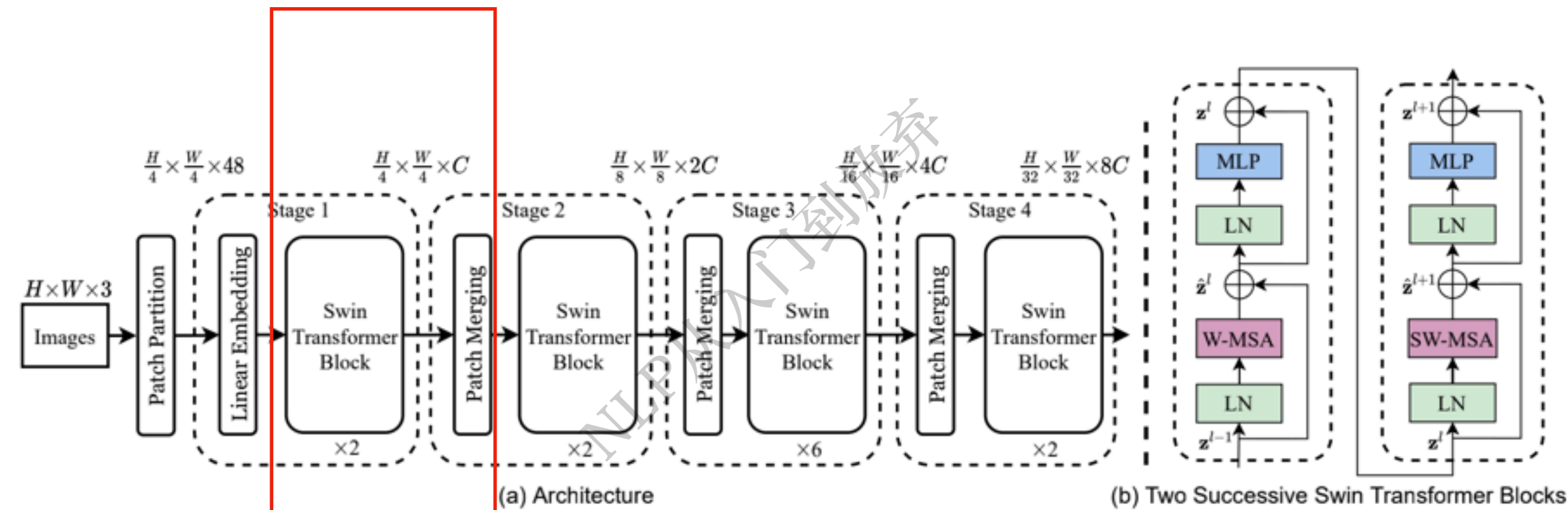
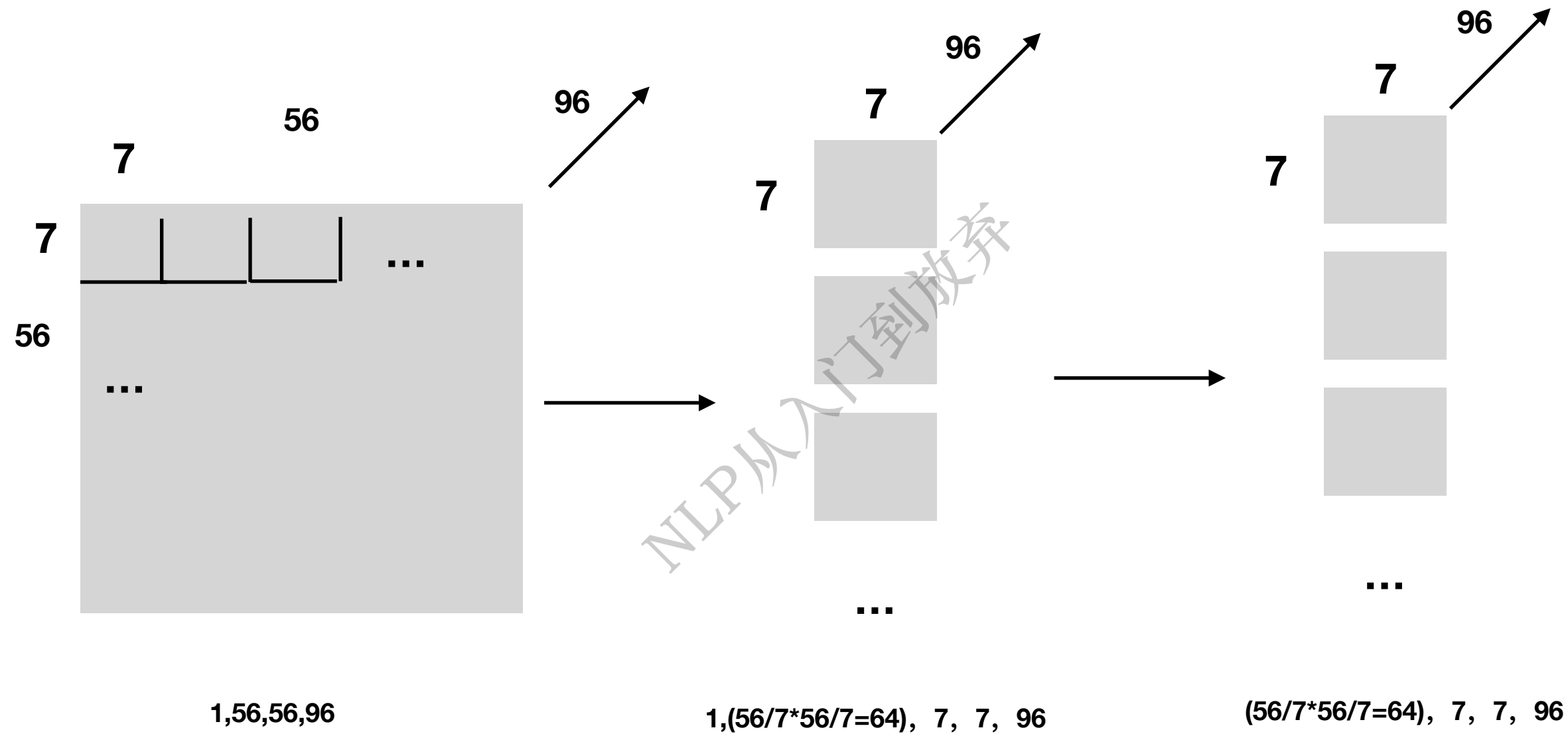
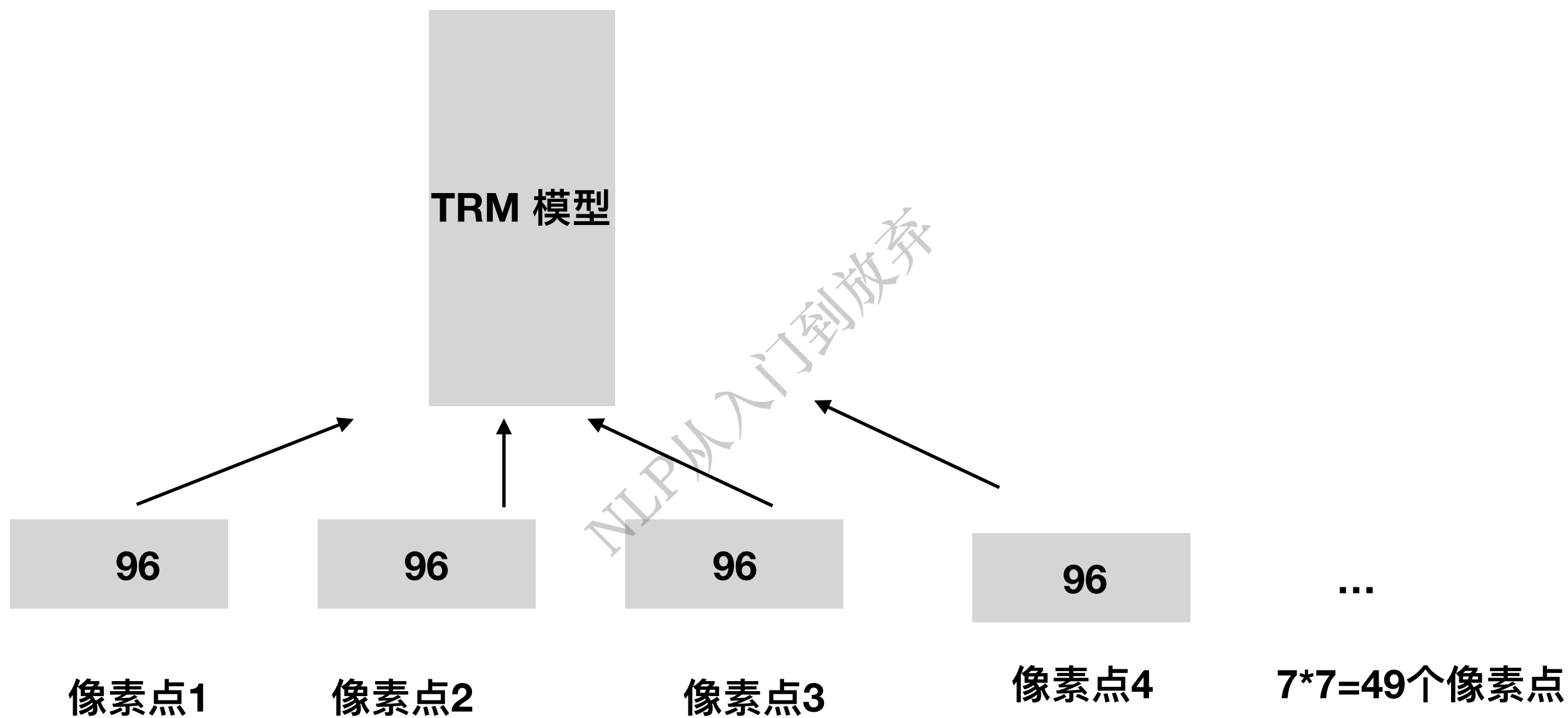


Figure 3. (a) The architecture of a Swin Transformer (Swin-T); (b) two successive Swin Transformer Blocks (notation presented with Eq. (3)). W-MSA and SW-MSA are multi-head self attention modules with regular and shifted windowing configurations, respectively.

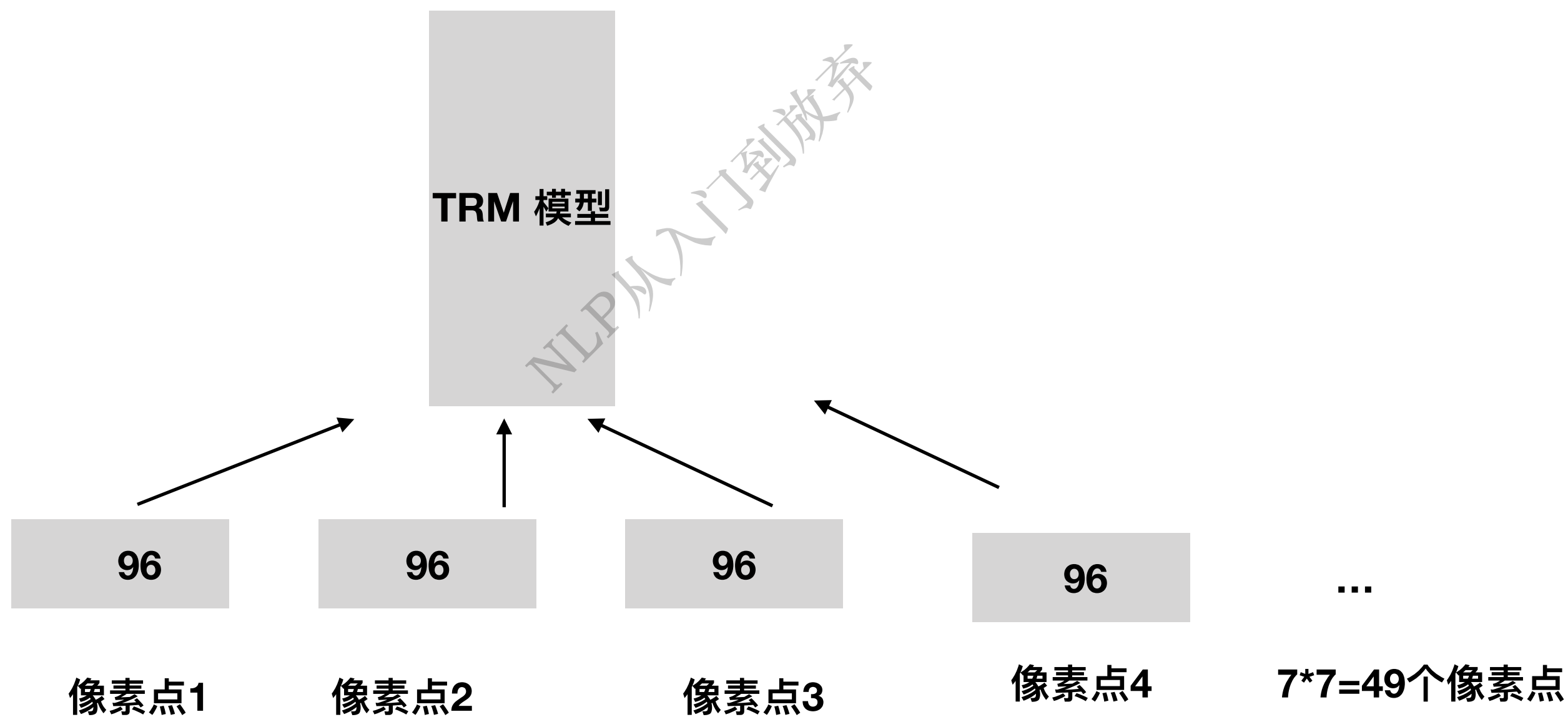
切分窗口：win_size=7一直保持不变



把一个图片的win数目移到bs维度上，就是整个patch有多少w



很简单：把每个元素点作为token，96个通道数作为token维度



位置编码问题

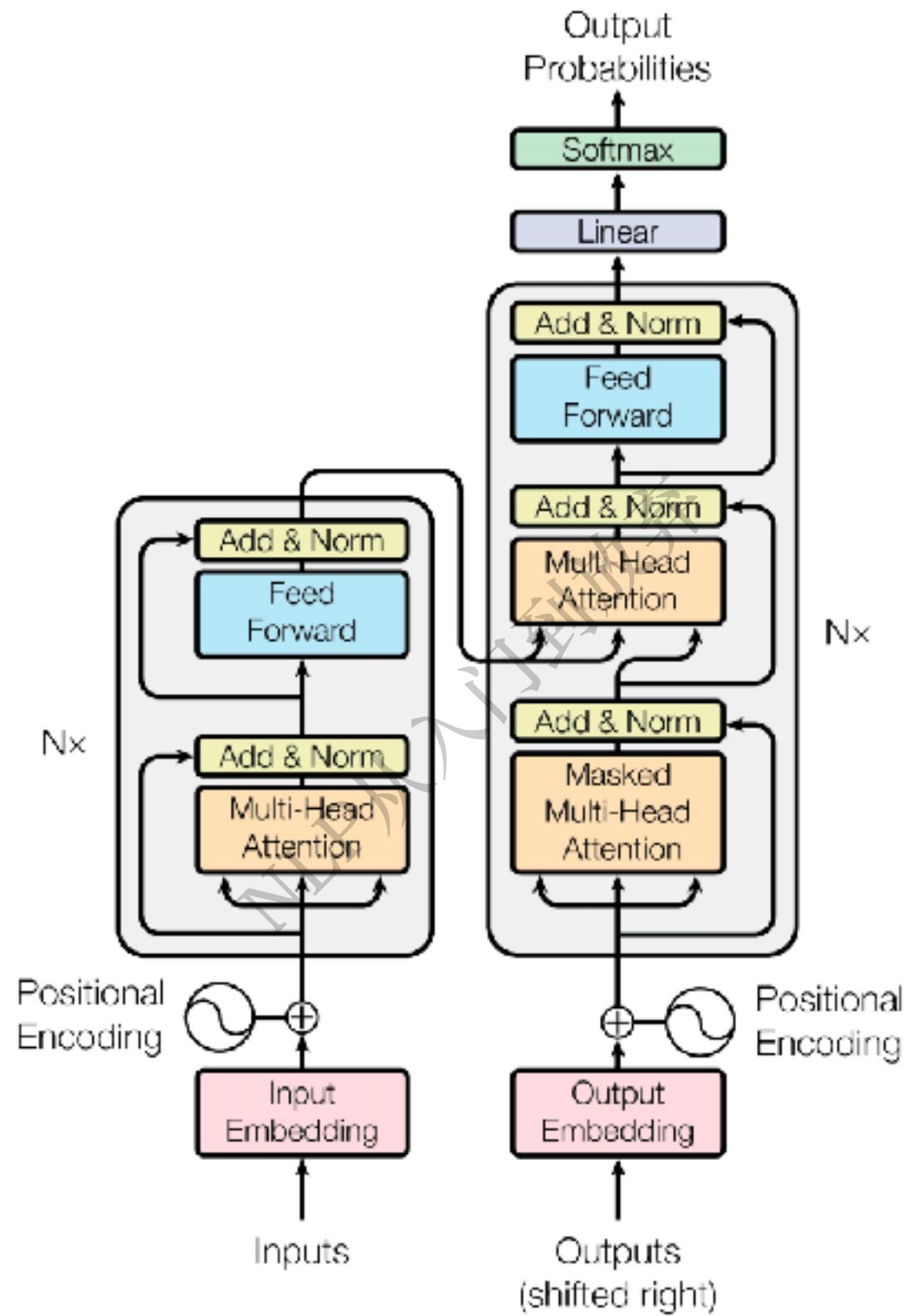
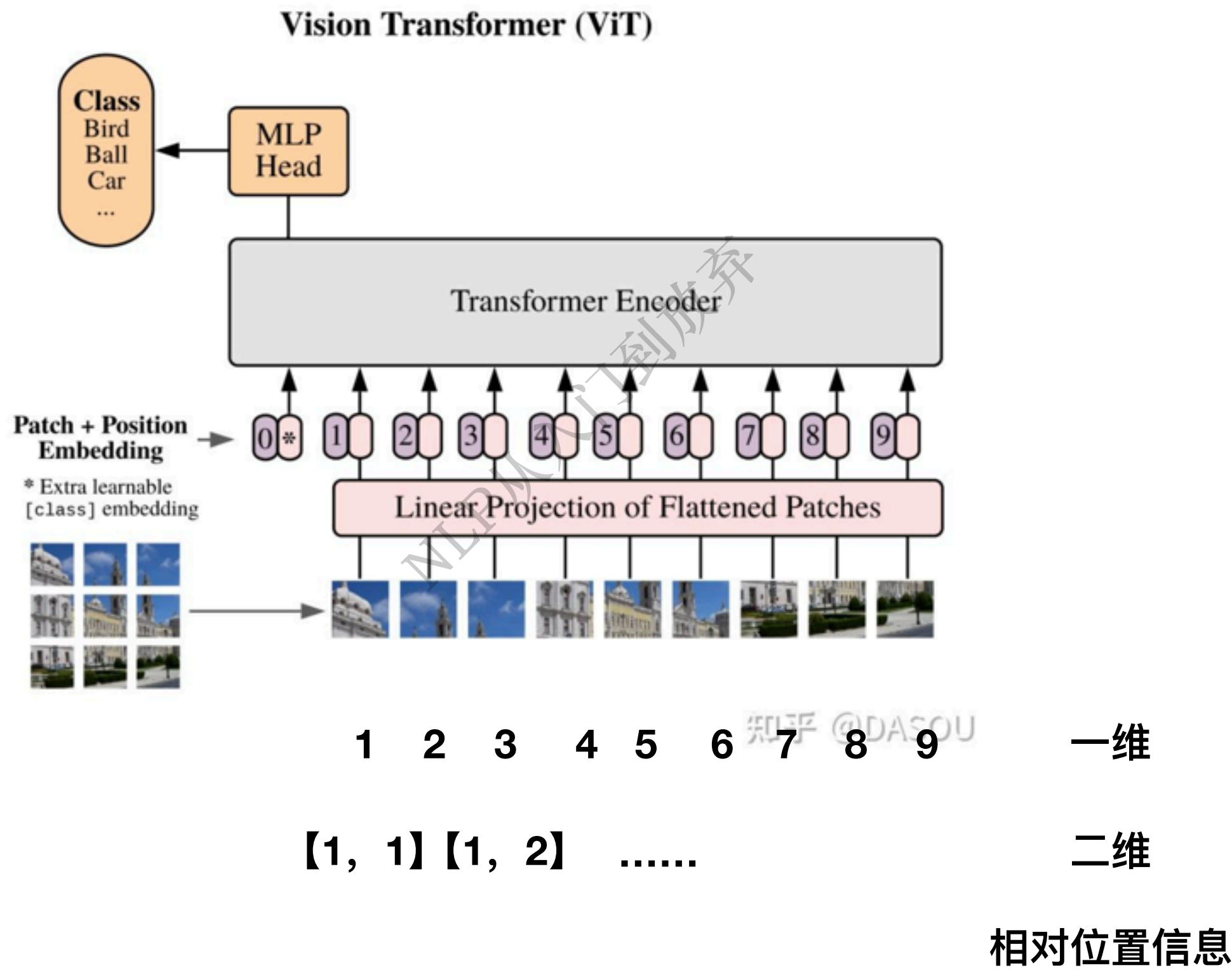


Figure 1: The Transformer - model architecture.

VIT中的位置编码：可学习的参数



SwinTRM位置编码有两点不同：

- 1. 加的位置不同：放在了att矩阵中**
- 2. 使用的是相对位置信息而不是绝对位置信息**

NLP从入门到放弃

最重要的相对位置编码讲解：

$$\text{Attention}(Q, K, V) = \text{SoftMax}(QK^T / \sqrt{d} + B)V, \quad (4)$$

QK转置相乘除以根号dk，然后加上相对位置编码

所以你需要知道**QK转置相乘除以根号dk**的输出形状是啥，你猜知道你的**B**是个啥形状

思考一下？

QK转置相乘除以根号dk的输出形状是啥？

NLP从入门到放弃

softmax函数里面这个东西本质是计算每个字符对每个字符的相似性对吧

形状不就应该下面是这种吗？

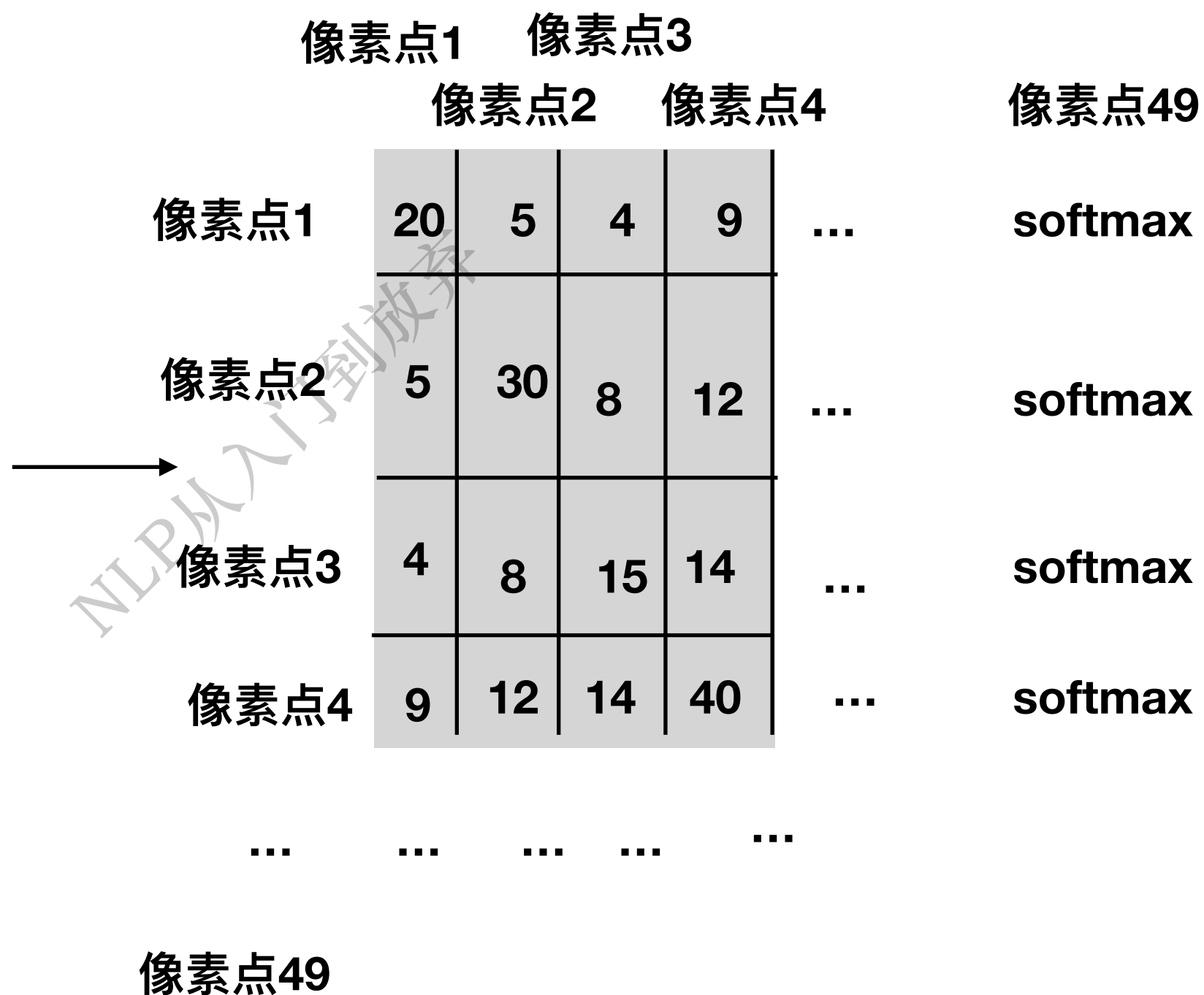
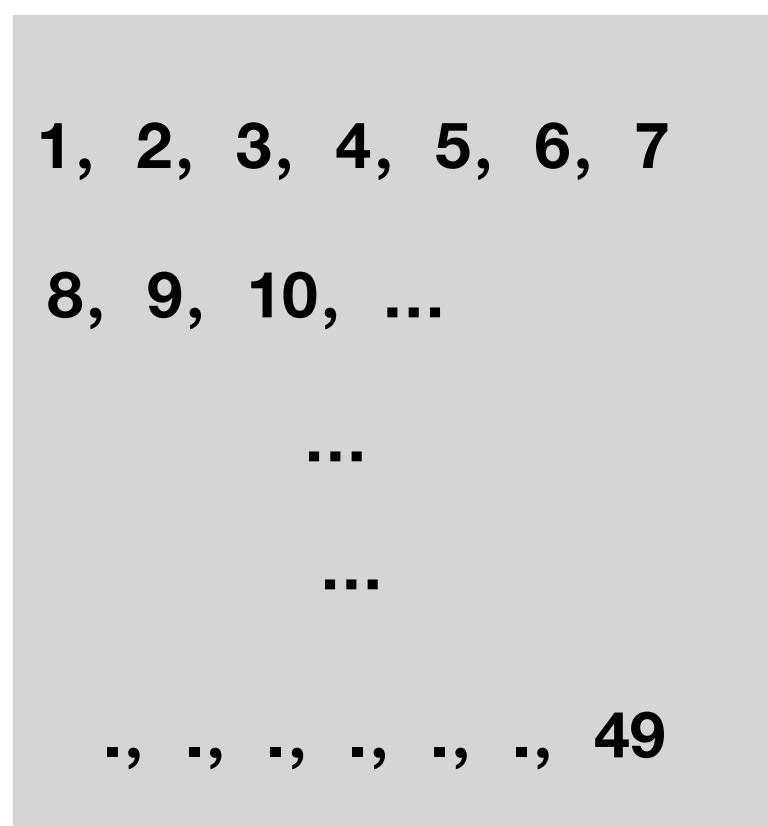
Seq_len*seq_len

只有这种形状，才能遍历到每个单词对每个单词的形状

	卷	起	来	吧	
卷	20	5	4	9	softmax
起	5	30	8	12	softmax
来	4	8	15	14	softmax
吧	9	12	14	40	softmax

一个7*7的图像，每个像素点对每个像素点的相似性，形状不就应该seq_le*seq_le

也就是49*49



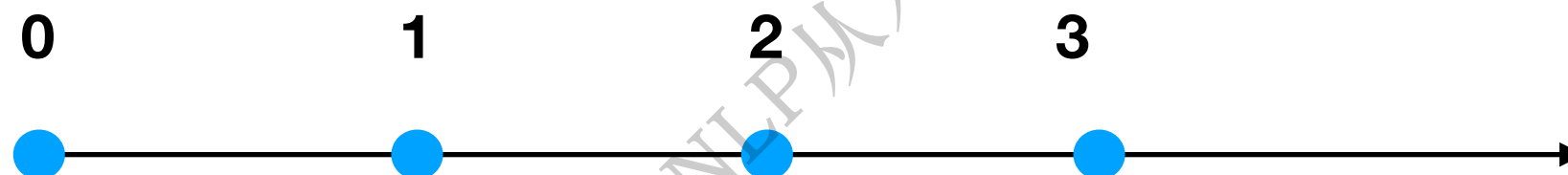
所以我的相对位置编码也必须是一个 49×49 的

7×7 窗口对应的这个 49×49 太大了，我以 $\text{win_size}=2$ 为例

$\text{win_size}=2$ ，attention矩阵就是： 4×4 对吧，在我接受范围之内啊

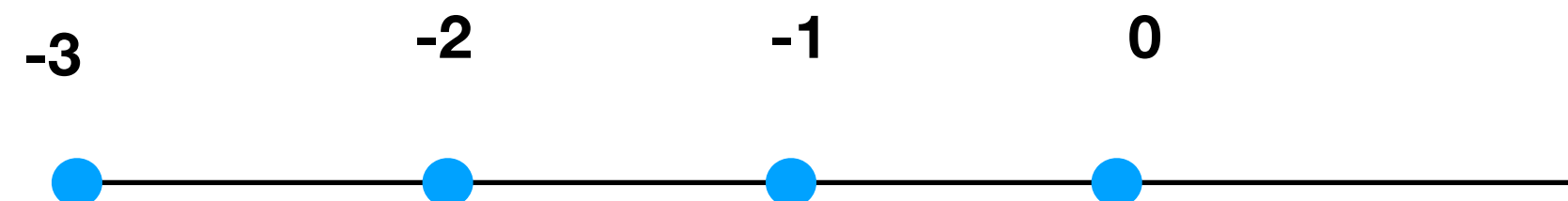
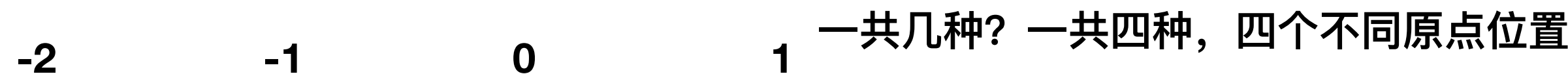
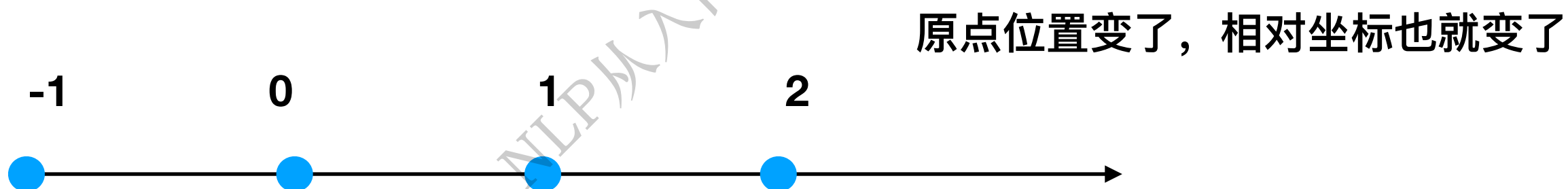
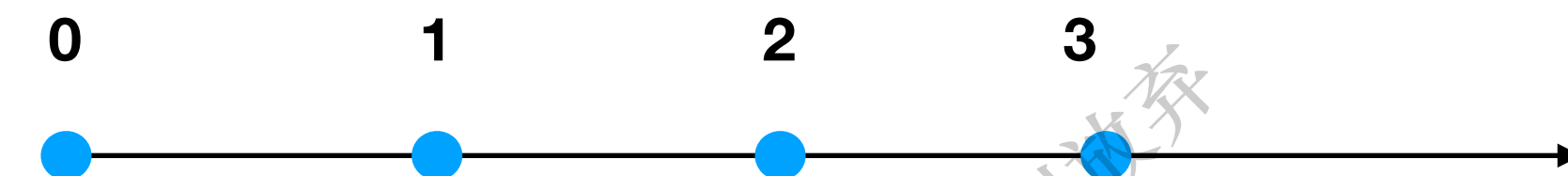
什么是绝对位置编码：一种绝对位置信息

绝对位置信息有很多种，就用最简单的例子来说啊



我的锚点其实没动，我的原点这种位置没变啊

什么是相对位置信息



刚才是一个线段的绝对位置信息和相对位置，再看网格绝对

网格绝对位置

像素1	像素2
像素3	像素4

0	1
2	3

网格相对位置信息

像素1	像素2
像素3	像素4

0	1
2	3

-1	0
1	2

-2	-1
0	1

-3	-2
-1	0

怎么把四种相对位置信息融合起来，放入到我attention矩阵呢？

像素1 像素2 像素3 像素4

像素1	20	5	4	9
像素2	5	30	8	12
像素3	4	8	15	14
像素4	9	12	14	40

0	1
2	3

-1	0
1	2

-2	-1
0	1

-3	-2
-1	0

像素1 像素2 像素3 像素4

像素1	20	5	4	9	以1为原点
像素2	5	30	8	12	以2为原点
像素3	4	8	15	14	以3为原点
像素4	9	12	14	40	以4为原点

像素1 像素2 像素3 像素4

像素1	20	5	4	9
像素2	5	30	8	12
像素3	4	8	15	14
像素4	9	12	14	40

以1为原点

0	1
2	3

以2为原点

-1	0
1	2

以3为原点

-2	-1
0	1

以4为原点

-3	-2
-1	0

0	1	2	3
---	---	---	---

-1	0	1	2
----	---	---	---

-2	-1	0	1
----	----	---	---

-3	-2	-1	0
----	----	----	---

SwinTRM中的相对位置编码是怎么搞的呢？很简单，就是把我刚才讲的
位置信息全部变为二维

像素1	像素2
像素3	像素4

0, 0	0, 1
1, 0	1, 1

0, -1	0, 0
1, -1	1, 0

-1, 0	-1, 1
0, 0	0, 1

-1, -1	-1, 0
0, -1	0, 0

0, 0	0, 1
1, 0	1, 1

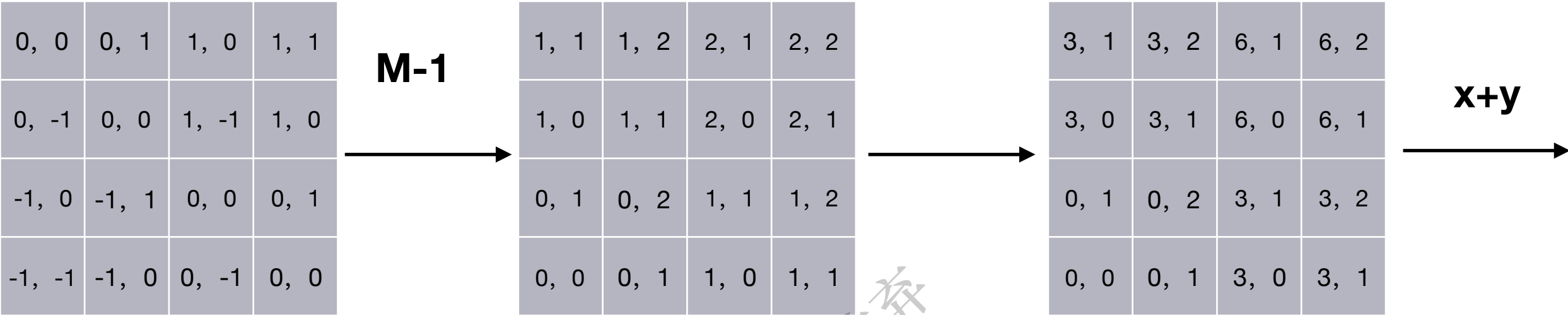
0, -1	0, 0
1, -1	1, 0

-1, 0	-1, 1
0, 0	0, 1

-1, -1	-1, 0
0, -1	0, 0

0, 0	0, 1	1, 0	1, 1
0, -1	0, 0	1, -1	1, 0
-1, 0	-1, 1	0, 0	0, 1
-1, -1	-1, 0	0, -1	0, 0

0维度*(2M-1)



索引|对应到参数

4	5	7	8
3	4	6	7
1	2	4	5
0	1	3	4

+

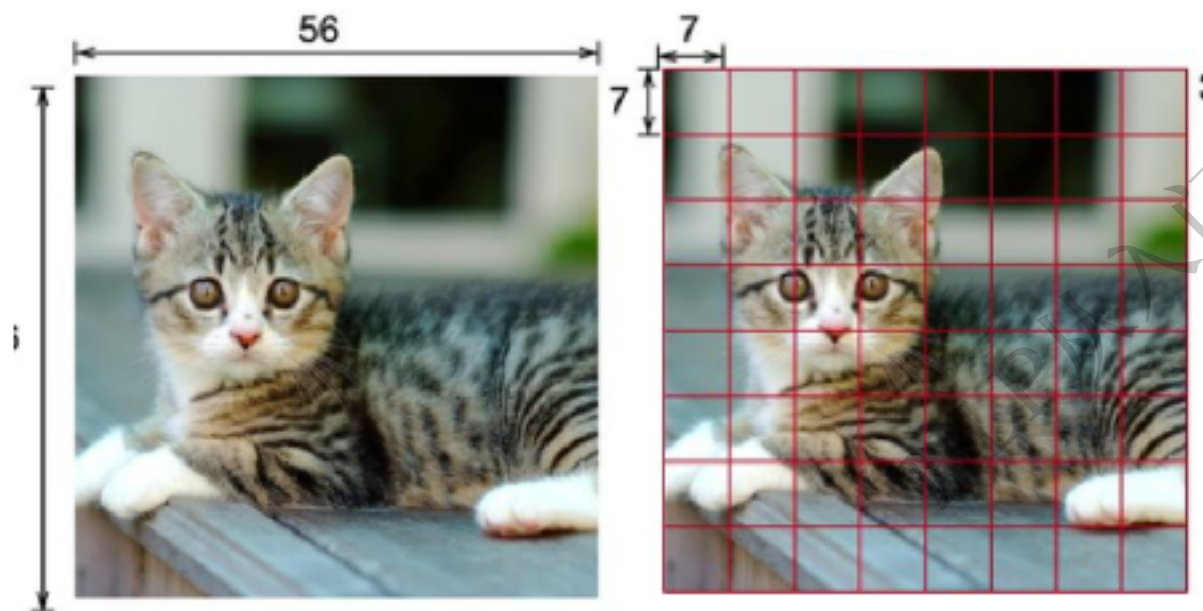
(2M-1)*(2M-1)

0.2	0.3	0.7	0.1	0.9	0.31	0.74	0.15
-----	-----	-----	-----	-----	------	------	------

像素1 像素2 像素3 像素4

像素1	20	5	4	9
像素2	5	30	8	12
像素3	4	8	15	14
像素4	9	12	14	40

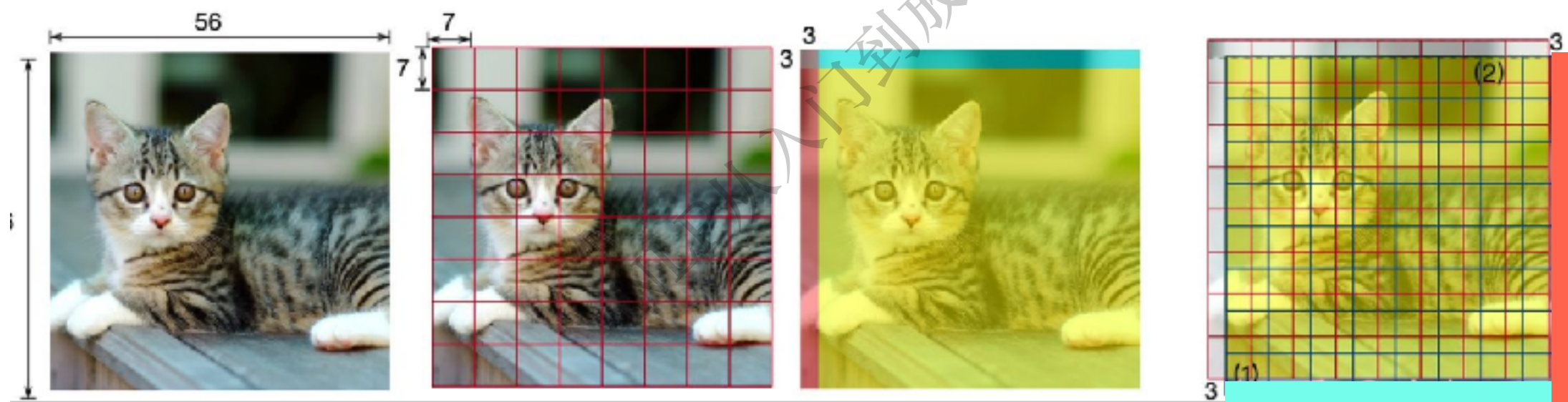
窗口注意力机制



7*7的win直接输入TRM就可以

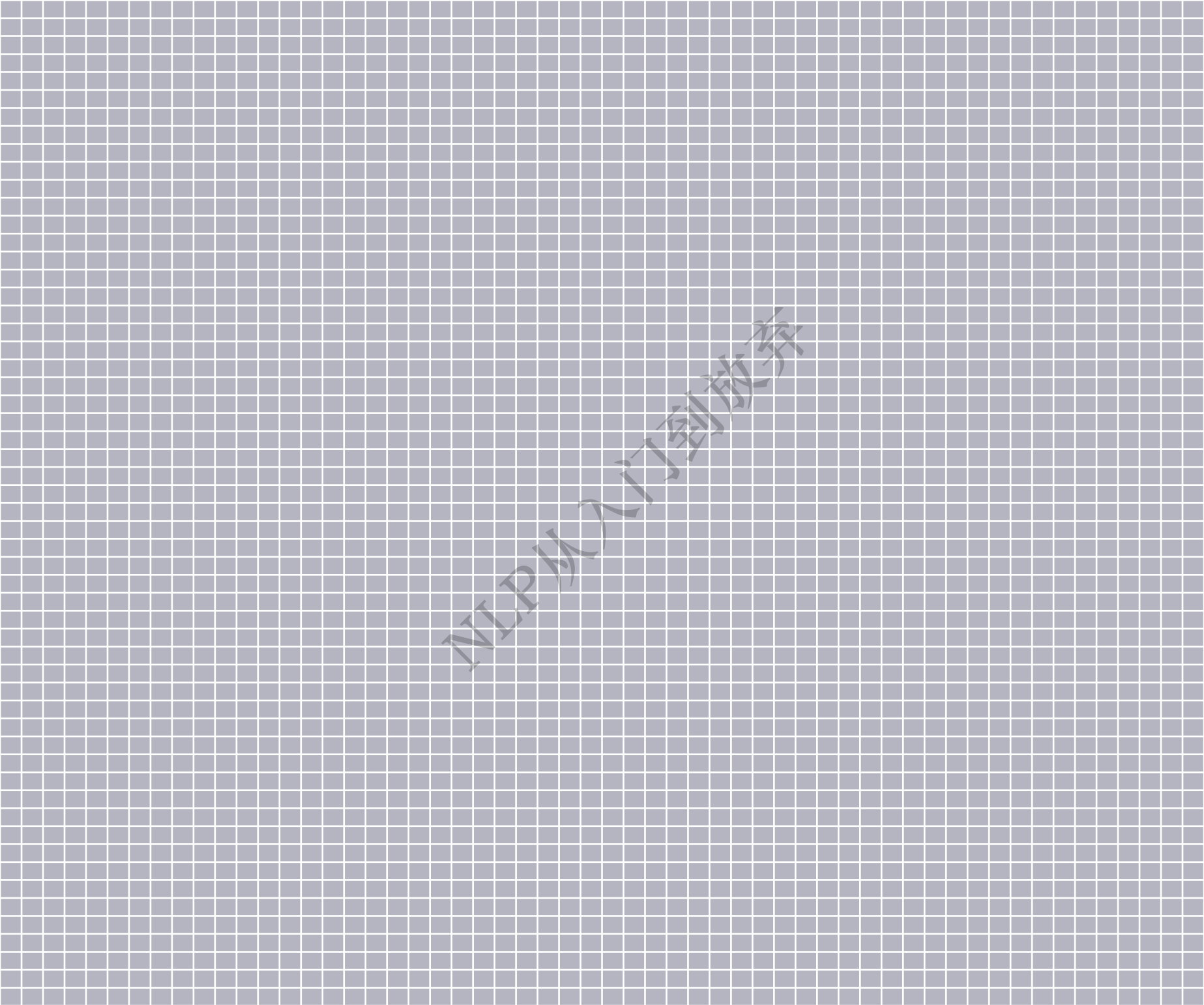
存在的一问题：窗口之间没有交互

但是我是这么理解的，没有交互只能说是在当前视野没有，
下一个阶段其实干感受野变大了，现在的窗口之间其实有交互，
但是在下一个阶段的窗口之间有没有交互了；所以为了在当前阶段就有交互，做了一个移动窗口注意力



但是问题是计算窗口att的时候，有的本不相邻的也会被计算：mask的重要性就出来了

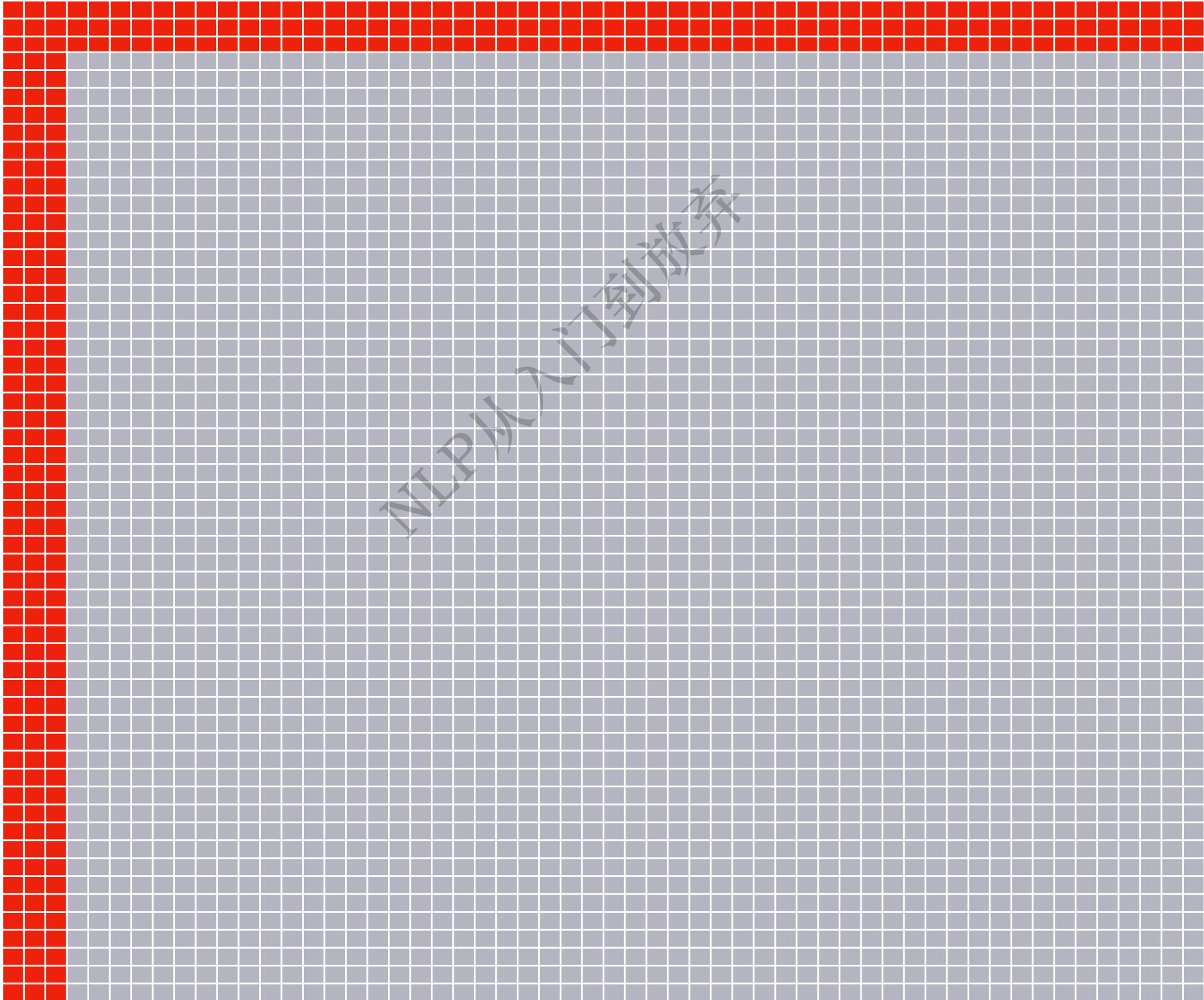
举个例子，我们从头说

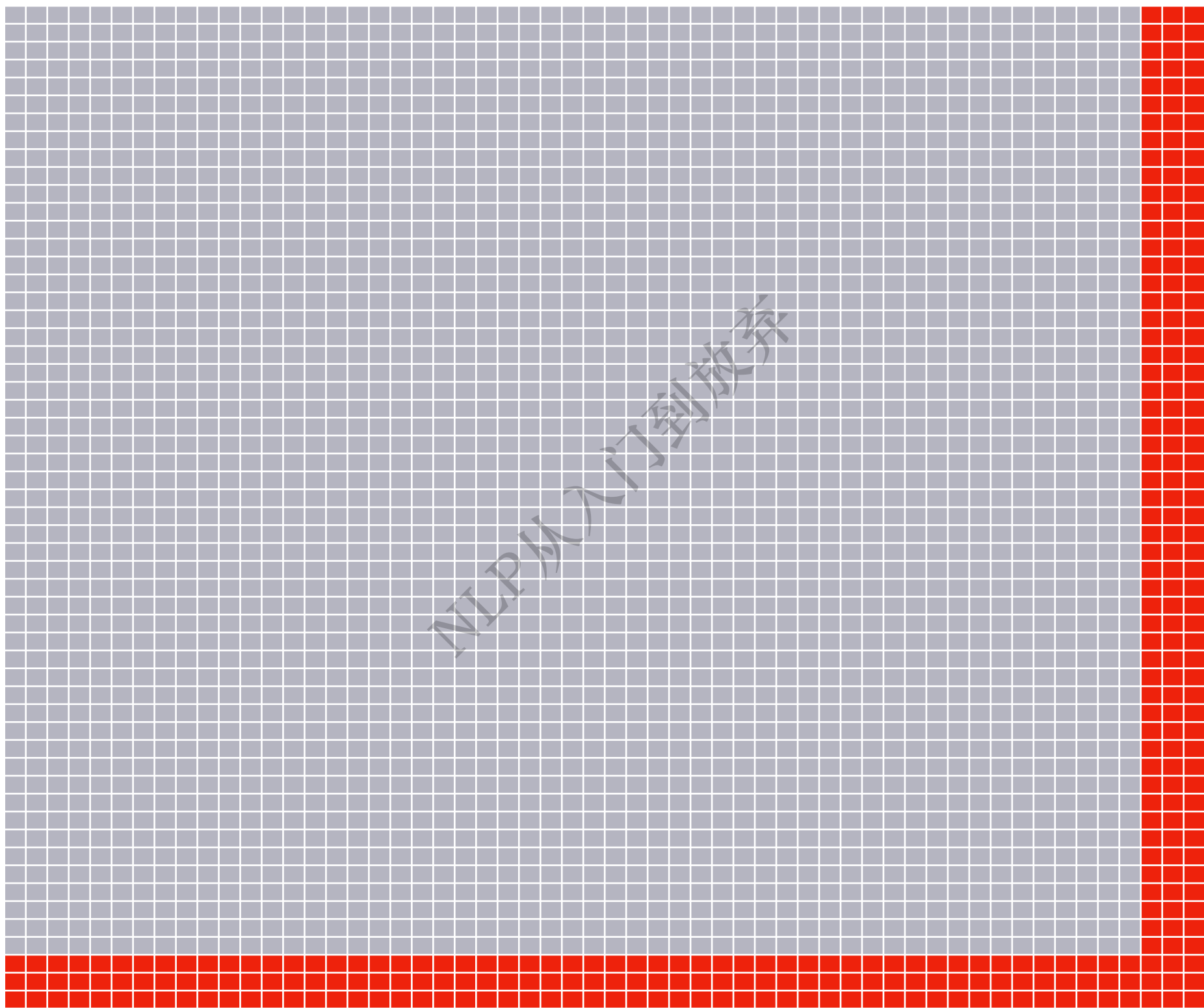


三行元素移动过去

56

56





有相邻的，有不相邻的，我们编上窗口索引

56

56-7=49

4

3

56-7=49

0

1

2

56

NLP从入门到放弃

4

3

4

5

3

6

7

8

56

56-7=49

4

3

7

7

56-7=49

0

NLP从入门到放弃

1

2

4

3

4

5

3

6

7

8

7

7

4	4	4	4	2	2	2
4	4	4	4	2	2	2
4	4	4	4	2	2	2
4	4	4	4	2	2	2
6	6	6	6	8	8	8
6	6	6	6	8	8	8
6	6	6	6	8	8	8

NLP从入门到放弃

4	2
6	8

4	2	6	8
---	---	---	---

4
2
6
8

0	-2	2	4
2	0	4	6
-2	-4	0	2
-4	-6	-2	0

4	2	6	8
---	---	---	---

4	0	-100	-100	-100
2	-100	0	-100	-100
6	-100	-100	0	-100
8	-100	-100	-100	0

4	2	6	8
---	---	---	---

	4	2	6	8
4	20	5	4	9
2	5	30	8	12
6	4	8	15	14
8	9	12	14	40

4
2
6
8

0	-100	-100	-100
-100	0	-100	-100
-100	-100	0	-100
-100	-100	-100	0

7

7

4	4	4	4	2	2	2
4	4	4	4	2	2	2
4	4	4	4	2	2	2
4	4	4	4	2	2	2
6	6	6	6	8	8	8
6	6	6	6	8	8	8
6	6	6	6	8	8	8

4	4	4	4	2	2	2	4	4	4	4	2	2	2	4	4	4	4	2	2	2	4	4	4	4	2	2	2	6	6	6	6	8	8	8	6	6	6	6	8	8	8	6	6	6	6	8	8	8
---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---

两个tensor相减

两个tensor相减

那三行元素究竟是怎么转的

代码里对特征图移位是通过 `torch.roll` 来实现的，下面是示意图



如果是一个2乘以2 的窗口，问题在于8和5是没办法attention，他们没关系
但是8和12是在一个窗口的

复杂度的问题

$$\Omega(\text{MSA}) = 4hwC^2 + 2(hw)^2C, \quad (1)$$

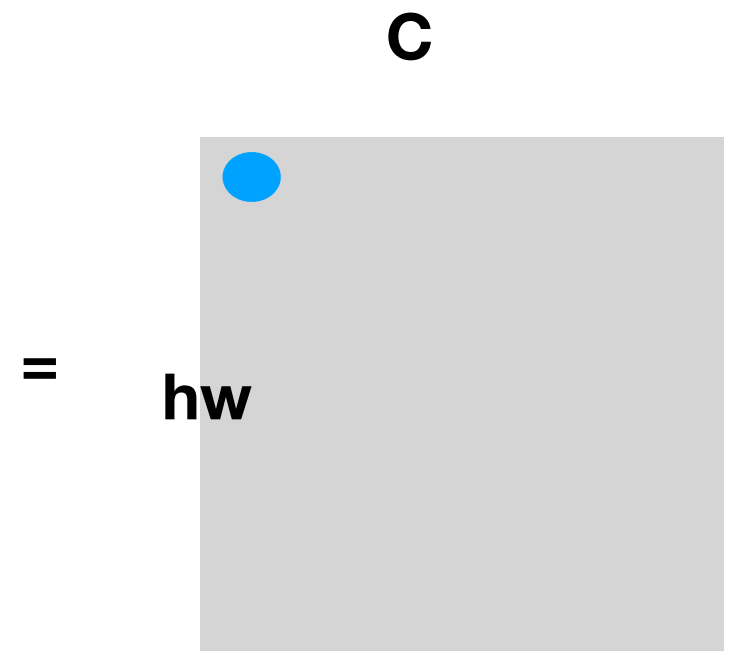
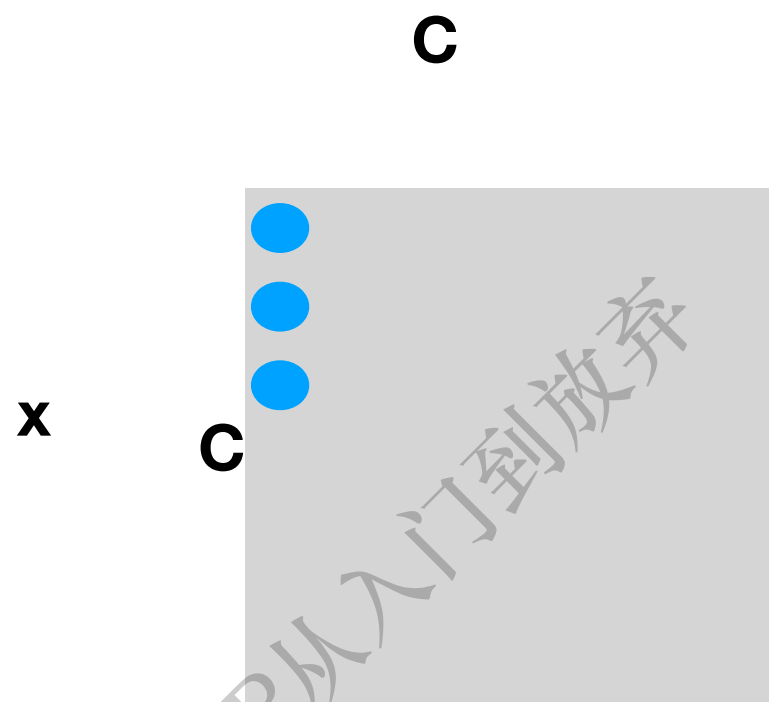
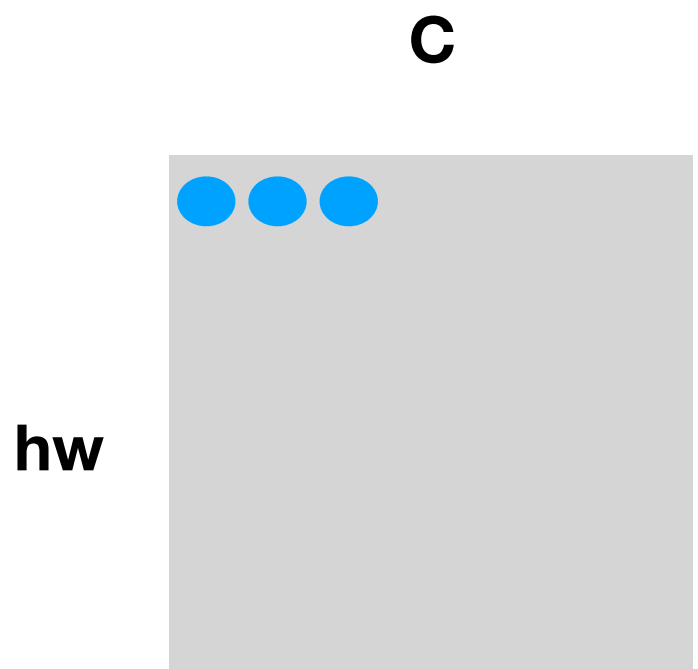
$$\Omega(\text{W-MSA}) = 4hwC^2 + 2M^2hwC, \quad (2)$$

$$\Omega(\text{MSA}) = 4hwC^2 + 2(hw)^2C, \quad (1)$$

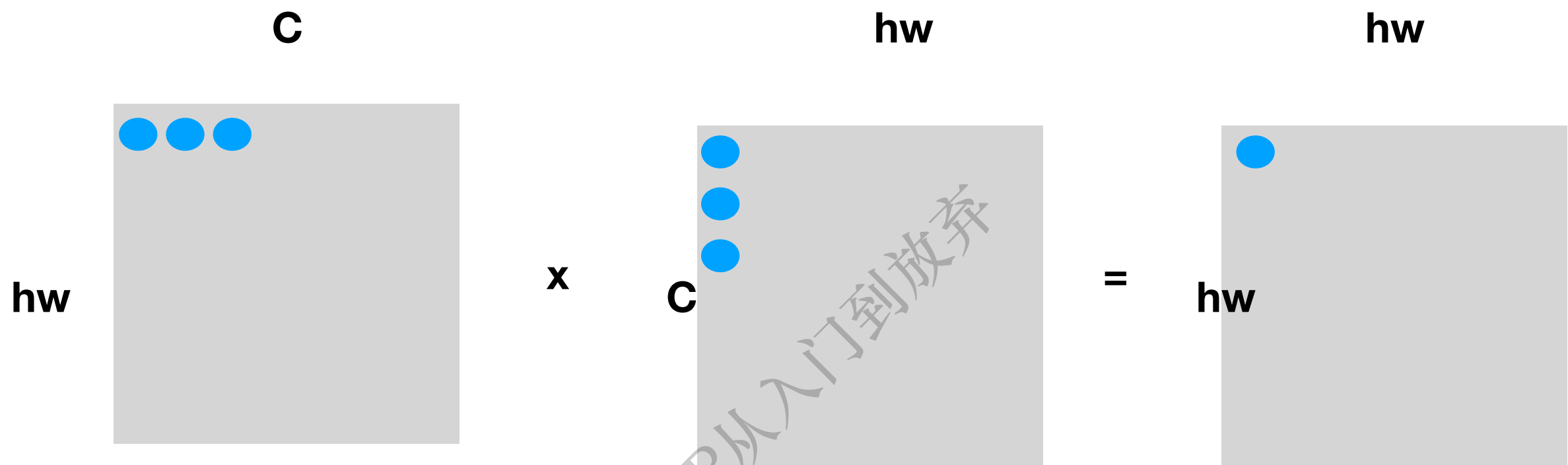
$$\Omega(\text{W-MSA}) = 4hwC^2 + 2M^2hwC, \quad (2)$$

1. 代码中的 `to_qkv()` 函数，即用于生成 Q, K, V 三个特征向量：其中 $Q = x \times W^Q, K = x \times W^K, V = x \times W^V$ 。 x 的维度是 (hw, C) ， W 的维度是 (C, C) ，那么这三项的复杂度是 $3hwC^2$ ；
2. 计算 QK^T ： Q, K, V 的维度均是 (hw, C) ，因此它的复杂度是 $(hw)^2C$ ；
3. softmax之后乘 V 得到 Z ：因为 QK^T 的维度是 (hw, hw) ，所以它的复杂度是 $(hw)^2C$ ；
4. Z 乘 W^Z 矩阵得到最终输出，对应代码中的 `to_out()` 函数：它的复杂度是 hwC^2 。

hw是长度，C是每个token的维度



hwC^2



v 矩阵

z

hw

c

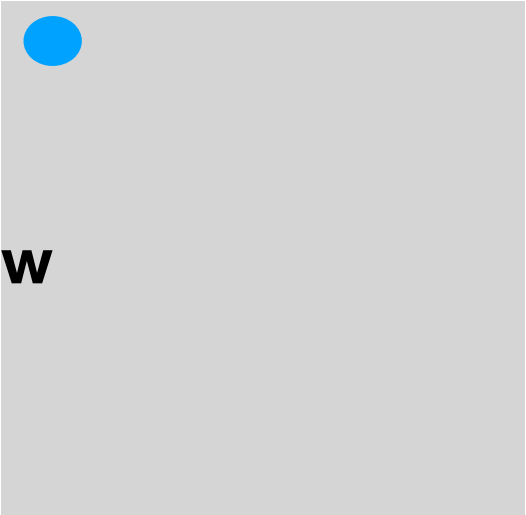
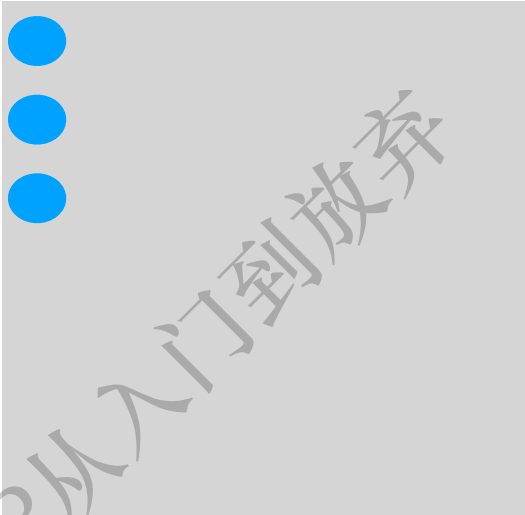
c

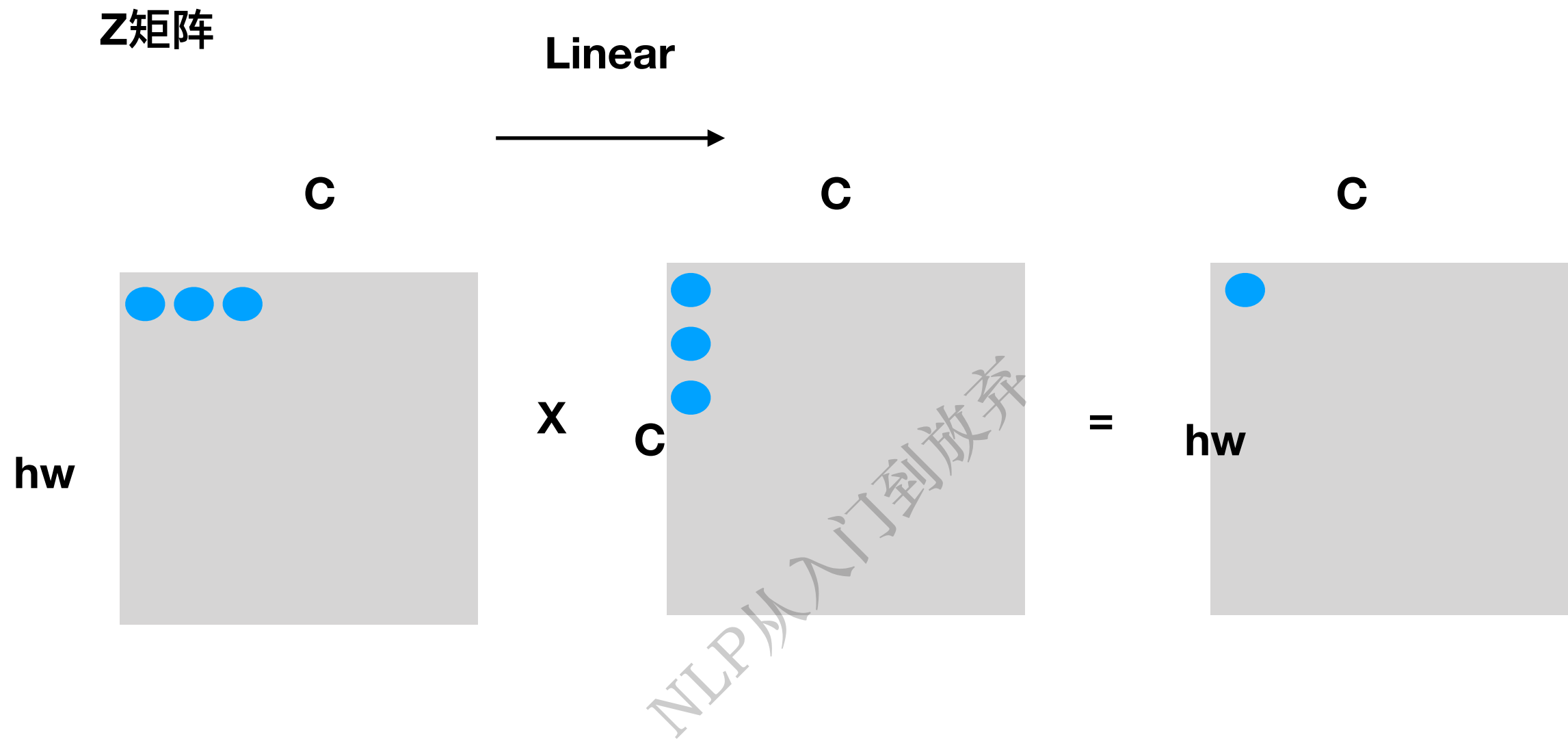
hw

x hw

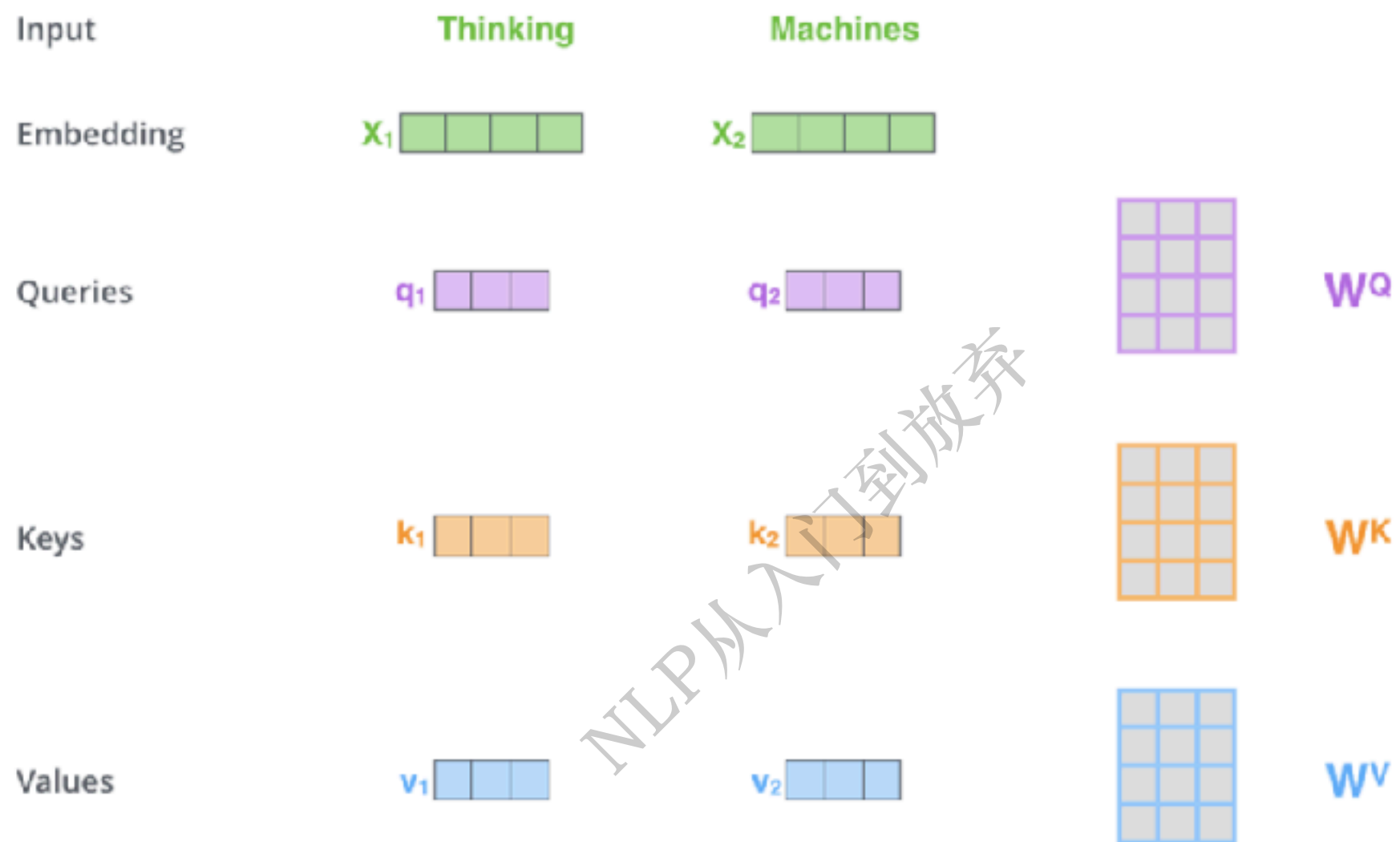
=

hw



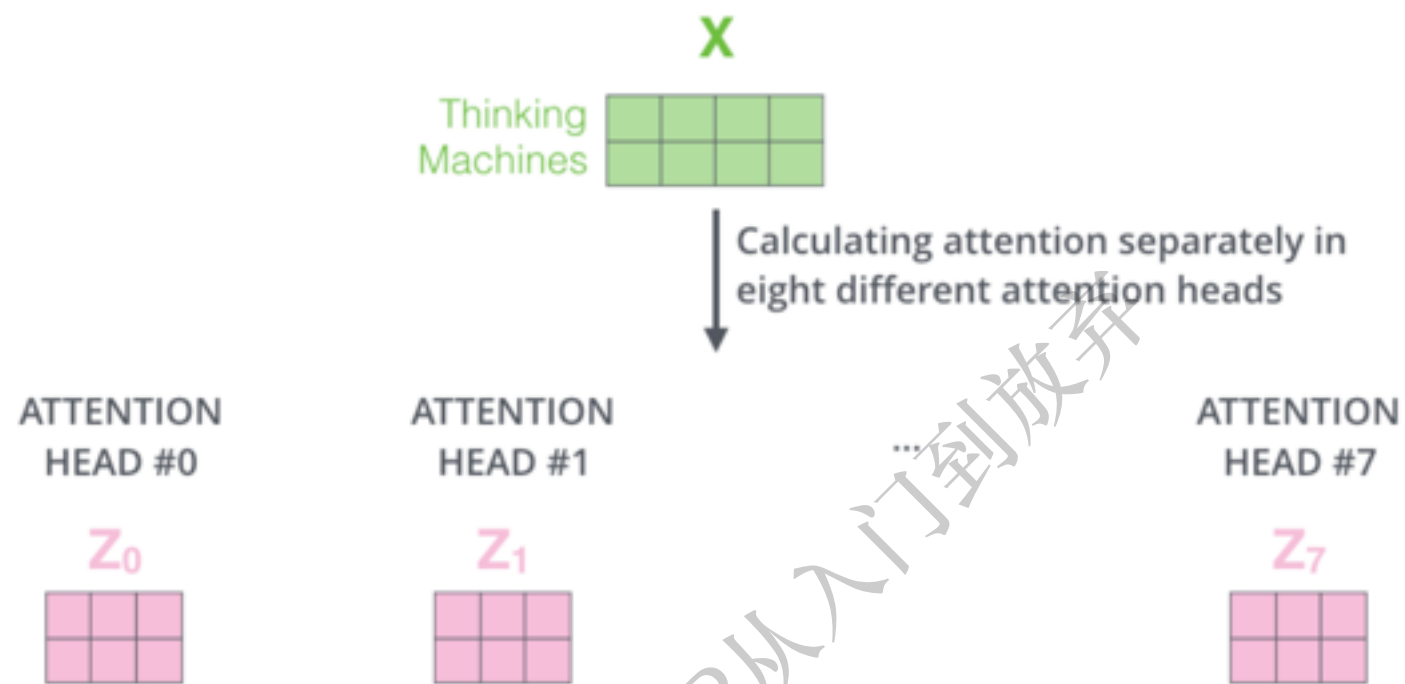


Z矩阵的最后的这个映射作用是什么？有的TRM代码实现可没这个步骤



一个头，四维度到了三维度

多头就是多个三维度拼在一起



TRM本身不改变形状的，不然不能堆叠，对从多个三拼接再到4维度就好了

窗口注意力机制

$$\Omega(\text{MSA}) = 4hwC^2 + 2(hw)^2C$$

—↑win

$$= 4 \boxed{\text{MM}} C^2 + 2(\boxed{\text{MM}})^2 C$$

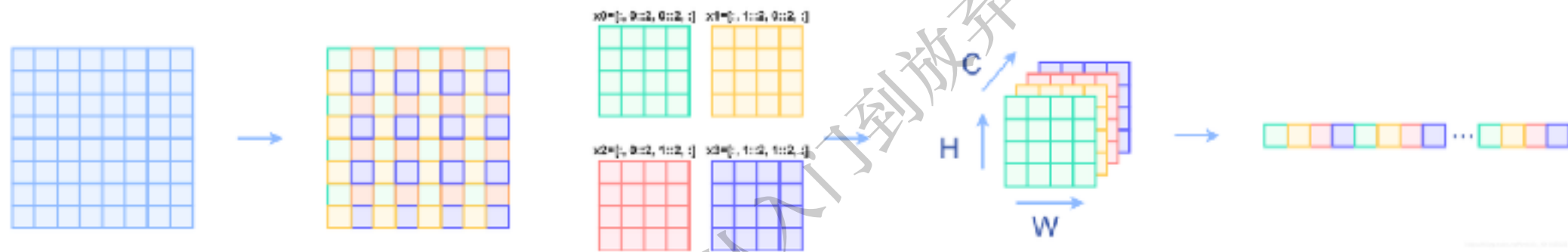
$h/M * w/M \uparrow \text{win}$

$(h/M)*(w/M)$

$$* 4 \boxed{\text{MM}} C^2 + 2(\boxed{\text{MM}})^2 C$$

$$\Omega(\text{W-MSA}) = 4hwC^2 + 2M^2hwC, \quad (2)$$

PatchMerging如何实现的：降采样



在行方向和列方向，间隔取2，通道维度会变成原先的4倍，接一个linear转为两倍

整体梳理一遍

H是224, W是224

